



TÉCNICO+
FORMAÇÃO AVANÇADA

Deep Learning

Gonçalo M. Correia

GALP 2023

Program for today

Generative models

Trust, interpretability, ethics

So far...

- All models discussed so far addressed **supervised learning**

So far...

- Example: Feed-forward neural networks
 - Data used:

$$D = \{(x_n, y_n), n = 1, \dots, N\}$$

where x_n is some vector of features, and y_n is the desired output (label, value)

- Network uses data to learn a correspondence between input space, \mathcal{X} , and output space, \mathcal{Y}

So far...

- Example: Convolutional neural networks

- Data used:

$$D = \{(x_n, y_n), n = 1, \dots, N\}$$

where x_n is some structured input (e.g., image, sound), and y_n is the desired output (label, value)

- Network uses data to learn a correspondence between \mathcal{X} and \mathcal{Y}

So far...

- Example: Recurrent neural networks

- Data used:

$$D = \{(x_n, y_n), n = 1, \dots, N\}$$

where x_n is a sequence of inputs (e.g., sentence), and y_n is the desired output (e.g., next word)

- Network uses data to learn a correspondence between \mathcal{X} and \mathcal{Y}

So far...

- Example: Sequence-to-sequence models
 - Data used:

$$D = \{(x_n, y_n), n = 1, \dots, N\}$$

where x_n is a sequence of inputs (e.g., sentence in a language), and y_n is a desired output sequence (e.g., sequence in a different language)

- Network uses data to learn a correspondence between \mathcal{X} and \mathcal{Y}

So far...

- Some of these models can be trained in a **self-supervised** manner
- The same data is used for input and for output
- Examples:
 - Auto-encoders
 - (Static) word embeddings
 - Large pre-trained models

So far...

- Some of these models can be trained in a self-supervised manner
- The same data is used for input and for output
- Self-supervised learning can be seen as a very particular form of **unsupervised learning**

Generative models

- What if we want to model the data itself?
- In other words, given a dataset

$$D = \{x_n, n = 1, \dots, N\}$$

we want to learn a model for $x_n, n = 1, \dots, N$.

Generative models

- A generative model is exactly a model of the data itself
- It is called “generative” since we can use it to generate synthetic data
- The model typically consists of learning a function

$$p_w : \mathcal{X} \rightarrow [0, 1]$$

so that

$$p_w(x) \approx \mathbb{P}[X = x]$$

Examples of generative models

- Auto-regressive networks
- Restricted Boltzmann machines
- Deep belief networks
- Deep Boltzmann machines
- Variational auto-encoders (VAEs)
- Generative adversarial networks (GANs)
- ...

Examples of generative models

- Auto-regressive networks
- Restricted Boltzmann machines
- Deep belief networks
- Deep Boltzmann machines
- Variational auto-encoders (VAEs)
- Generative adversarial networks (GANs)
- ...

Differentiable generator networks

- Both VAEs and GANs belong to a class of models known as **differentiable generator networks**
- These networks compute a conditional model $p_w(x \mid z)$
 - z is a (low-dimensional) latent variable
 - z captures different dimensions of variation of the data

How can we train p_w ?

Let's play the math a bit...

- From Bayes rule,

$$p_w(z \mid x) = \frac{p_w(x, z)}{p_w(x)}$$

or, equivalently,

$$p_w(x) = \frac{p_w(x, z)}{p_w(z \mid x)}$$

Let's play the math a bit...

- From Bayes rule,

$$p_w(z \mid x) = \frac{p_w(x, z)}{p_w(x)}$$

or, equivalently,

$$\log p_w(x) = \log p_w(x, z) - \log p_w(z \mid x)$$

Let's play the math a bit...

- Let's take an arbitrary distribution q over z ...

$$\log p_w(x) = \log p_w(x, z) - \log p_w(z \mid x)$$

Let's play the math a bit...

- Let's take an arbitrary distribution q over z ...

$$\log p_w(x) = \log p_w(x, z) - \log p_w(z \mid x) + \log q(z) - \log q(z)$$

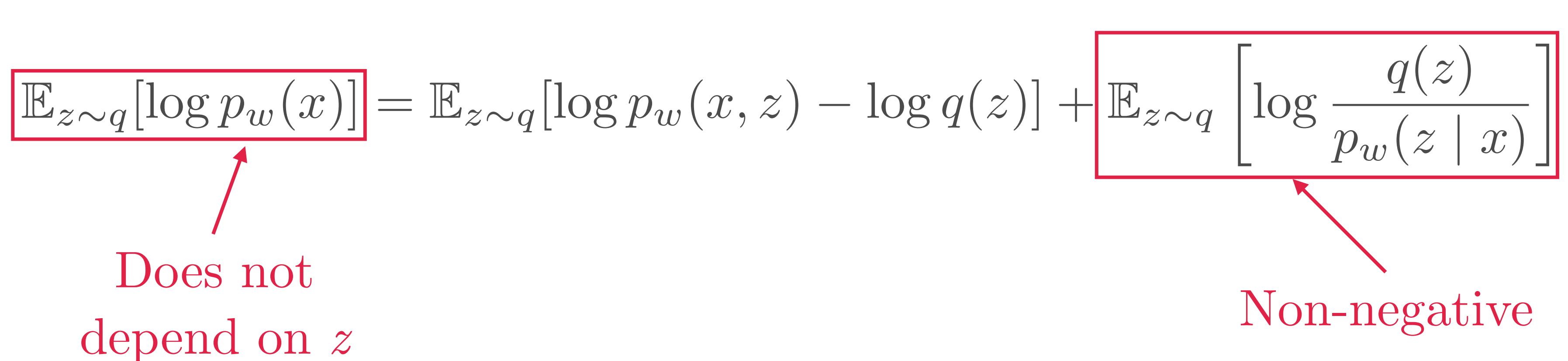
Let's play the math a bit...

- Let's take an arbitrary distribution q over z ...

$$\log p_w(x) = \log p_w(x, z) - \log q(z) + \log \frac{q(z)}{p_w(z | x)}$$

- Taking the expectation w.r.t. q ...

$$\mathbb{E}_{z \sim q}[\log p_w(x)] = \mathbb{E}_{z \sim q}[\log p_w(x, z) - \log q(z)] + \mathbb{E}_{z \sim q} \left[\log \frac{q(z)}{p_w(z | x)} \right]$$



Does not depend on z

Non-negative

Let's play the math a bit...

- Let's take an arbitrary distribution q over z ...

$$\log p_w(x) = \log p_w(x, z) - \log q(z) + \log \frac{q(z)}{p_w(z | x)}$$

- Taking the expectation w.r.t. q ...

$$\log p_w(x) \geq \mathbb{E}_{z \sim q} [\log p_w(x, z) - \log q(z)]$$

Evidence lower bound (ELBO)

ELBO

- Evidence lower bound (ELBO):

$$\log p_w(x) \geq \mathbb{E}_{z \sim q}[\log p_w(x, z) - \log q(z)]$$

- Idea:
 - If we push the right-hand side up, we push the left-hand side up
 - How?
 - Choosing q adequately

ELBO

- We can rewrite the ELBO as

$$\text{ELBO} = \boxed{\mathbb{E}_{z \sim q}[\log p_w(x \mid z)]} - \boxed{\text{KL}(q(z) \parallel p_{\text{prior}}(z))}$$

Maximize the ability to
recover x from z
(reconstruction)

Minimize the distance
between q and prior
(regularization)

What about q ?

- How can we select q ?
- Standard approach is to pick q to...
 - ... select z from x
 - ... be some “friendly distribution” (e.g., Gaussian)
 - ... be itself a neural network

Putting everything together...

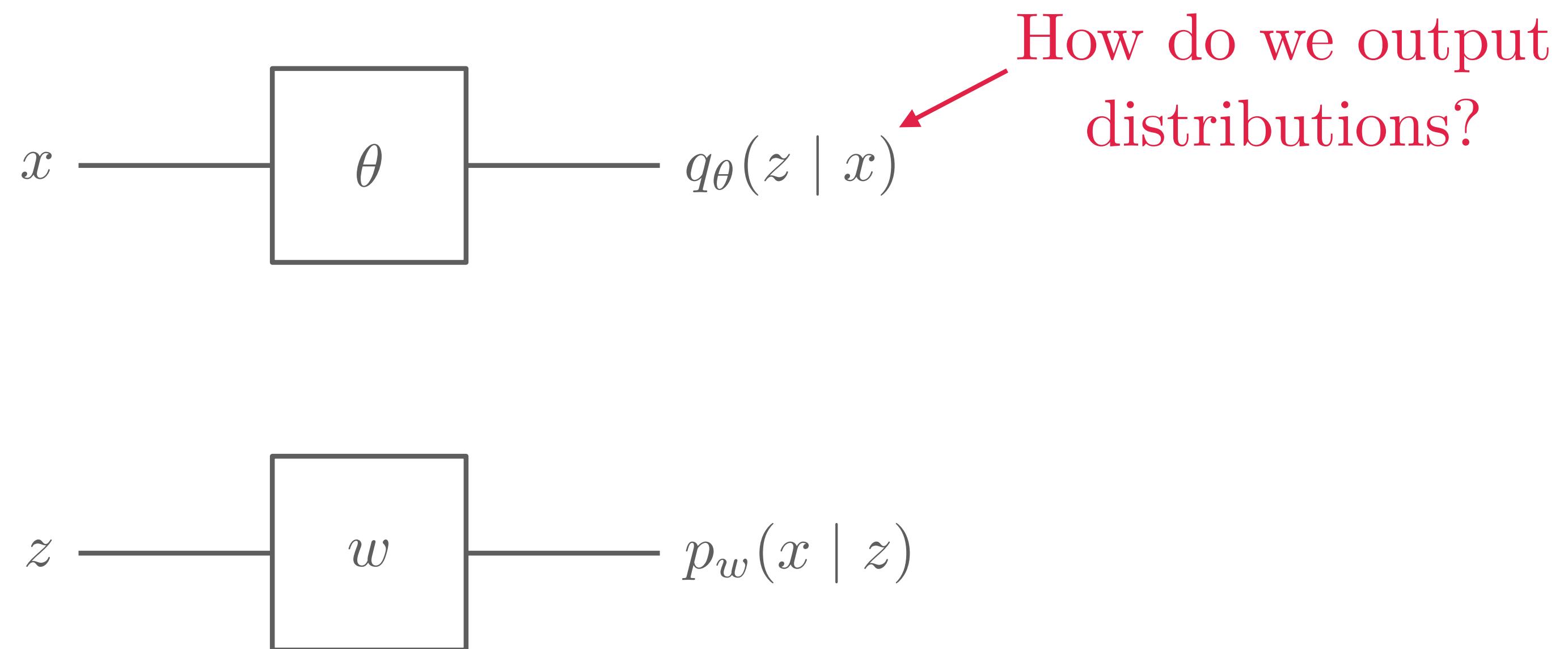
- Back to the ELBO:

$$\text{ELBO} = \mathbb{E}_{z \sim q} [\log p_w(x \mid z)] - \text{KL}(q_\theta(z \mid x) \parallel p_{\text{prior}}(z))$$

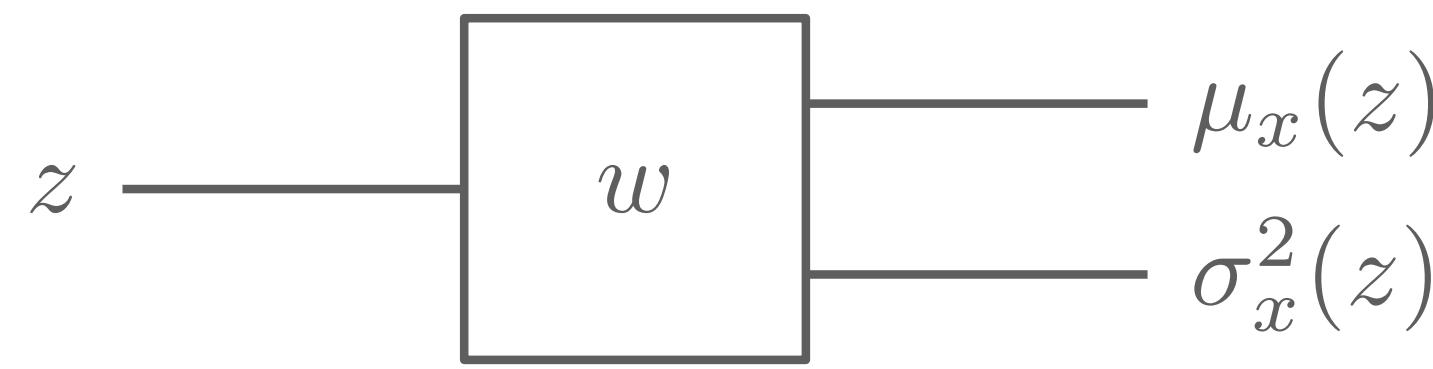
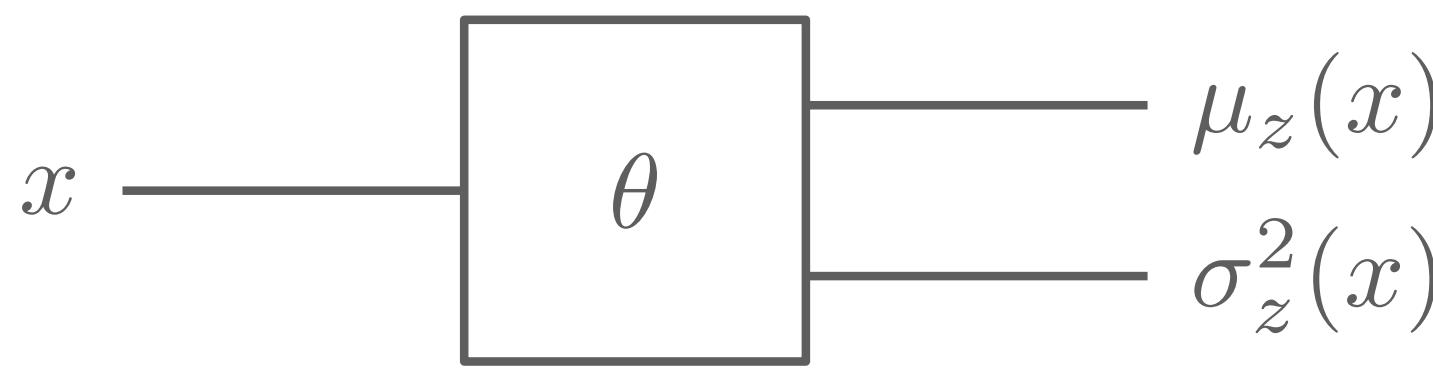
- We have 3 elements:

- A neural network that encodes $p_w(x \mid z)$
- A neural network that encodes $q_\theta(z \mid x)$
- The prior p_{prior} , usually a $\text{Normal}(0, I)$

Putting everything together...

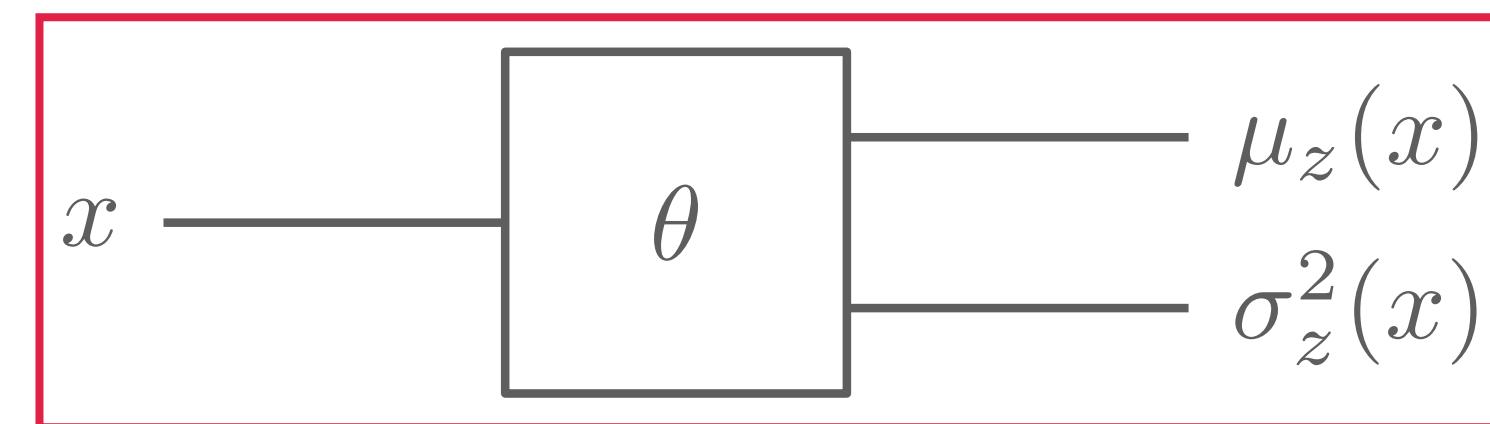


Putting everything together...

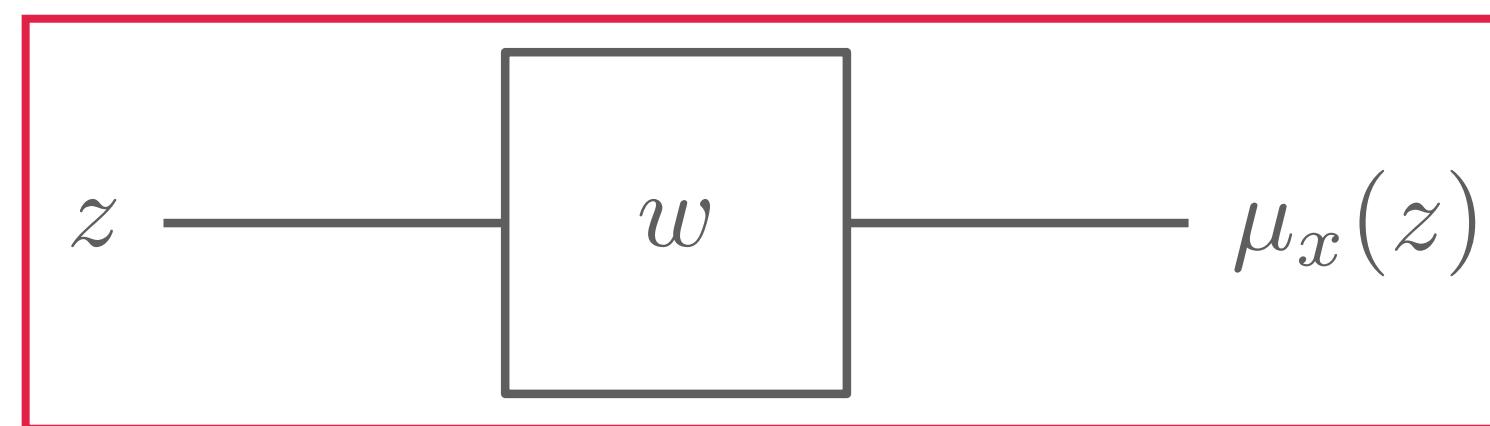


To make our life easier,
we usually consider
 $\sigma_x^2(z) = 1$

Putting everything together...

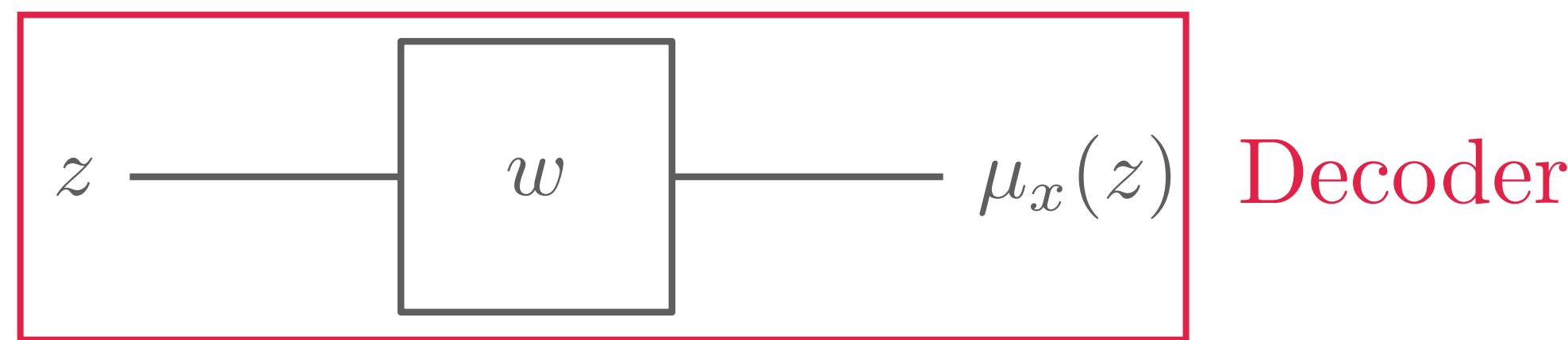
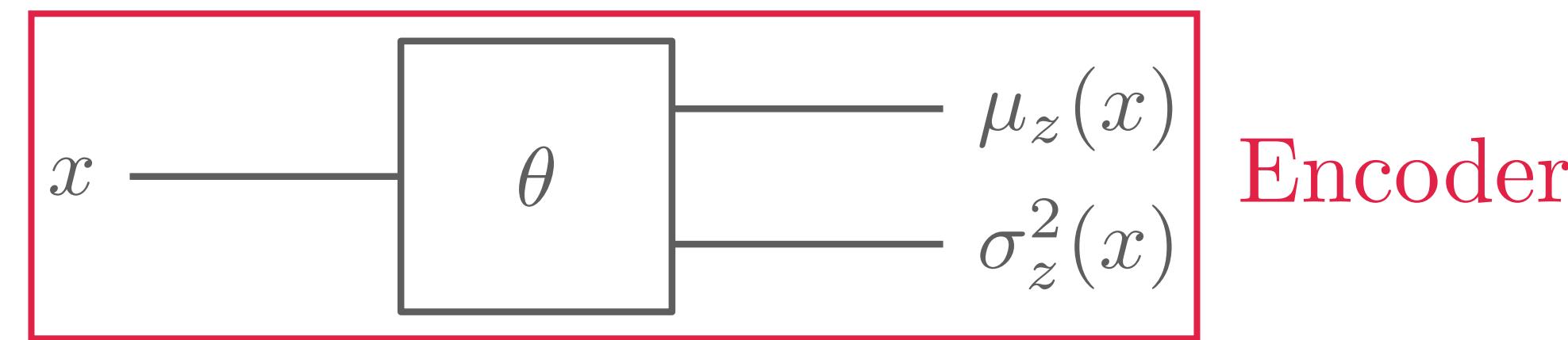


Returns a (sampled) representation z of x

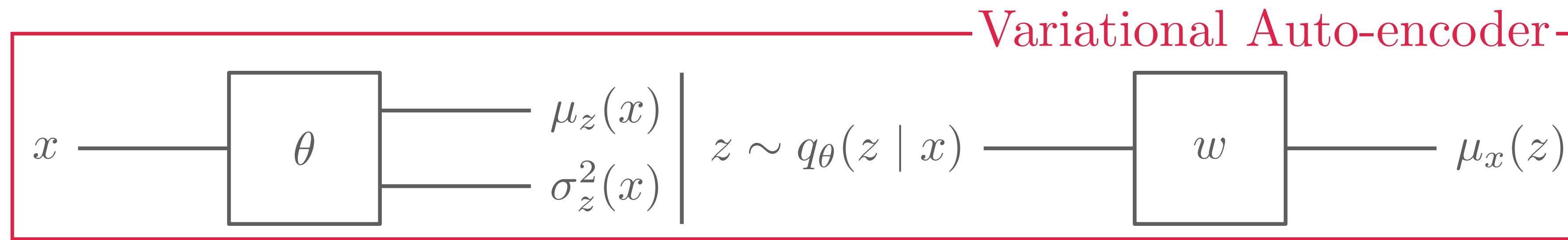


Recovers x from z

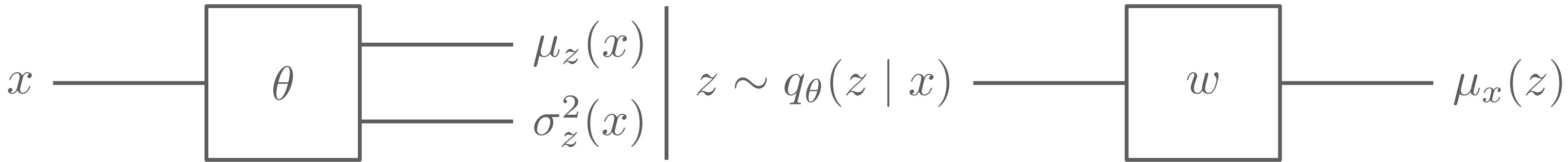
Putting everything together...



Putting everything together...



Training a VAE



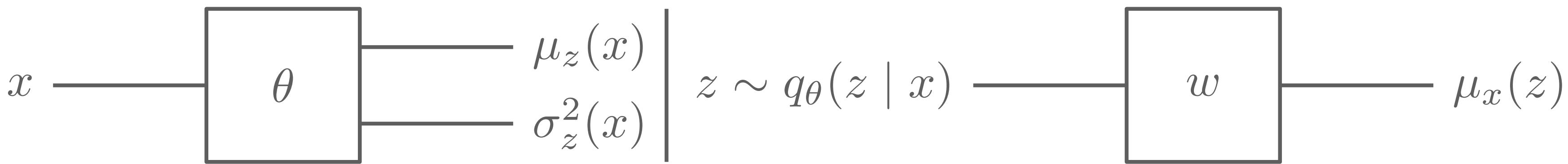
- We want to maximize the ELBO:

$$\text{ELBO} = \mathbb{E}_{z \sim q}[\log p_w(x \mid z)] - \text{KL}(q_\theta(z \mid x) \parallel p_{\text{prior}}(z))$$

or, equivalently, minimize the loss

$$L(w, \theta) = \text{KL}(q_\theta(z \mid x) \parallel p_{\text{prior}}(z)) - \mathbb{E}_{z \sim q}[\log p_w(x \mid z)]$$

Training a VAE



- We want to maximize the ELBO:

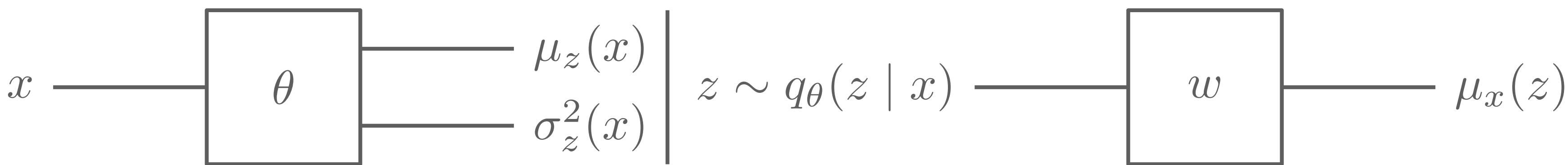
$$\text{ELBO} = \mathbb{E}_{z \sim q}[\log p_w(x | z)] - \text{KL}(q_\theta(z | x) \| p_{\text{prior}}(z))$$

or, equivalently, minimize the loss

$$L(w, \theta) = \mathbb{E}_{z \sim q}[\log q_\theta(z | x) - \log p_{\text{prior}}(z) - \log p_w(x | z)]$$

How do we deal
with the expectation?

Training a VAE



- We want to maximize the ELBO:

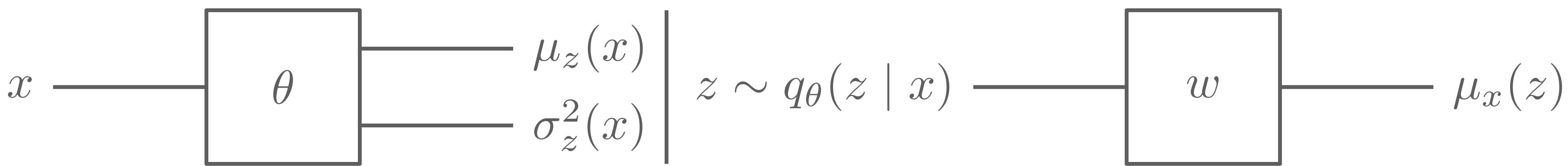
$$\text{ELBO} = \mathbb{E}_{z \sim q}[\log p_w(x \mid z)] - \text{KL}(q_\theta(z \mid x) \parallel p_{\text{prior}}(z))$$

or, equivalently, minimize the loss

$$L(w, \theta) = \frac{1}{M} \sum_{m=1}^M (\log q_\theta(z_m \mid x) - \log p_{\text{prior}}(z_m) - \log p_w(x \mid z_m))$$

Usually, $M = 1$

Training a VAE



- We want to maximize the ELBO:

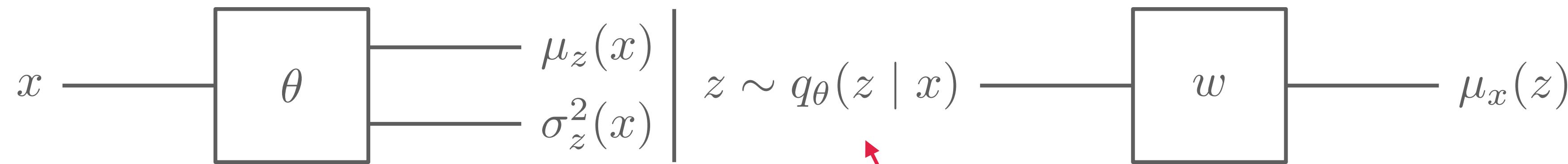
$$\text{ELBO} = \mathbb{E}_{z \sim q}[\log p_w(x | z)] - \text{KL}(q_\theta(z | x) \| p_{\text{prior}}(z))$$

or, equivalently, minimize the loss

$$L(w, \theta) = \log q_\theta(z | x) - \log p_{\text{prior}}(z) - \log p_w(x | z)$$

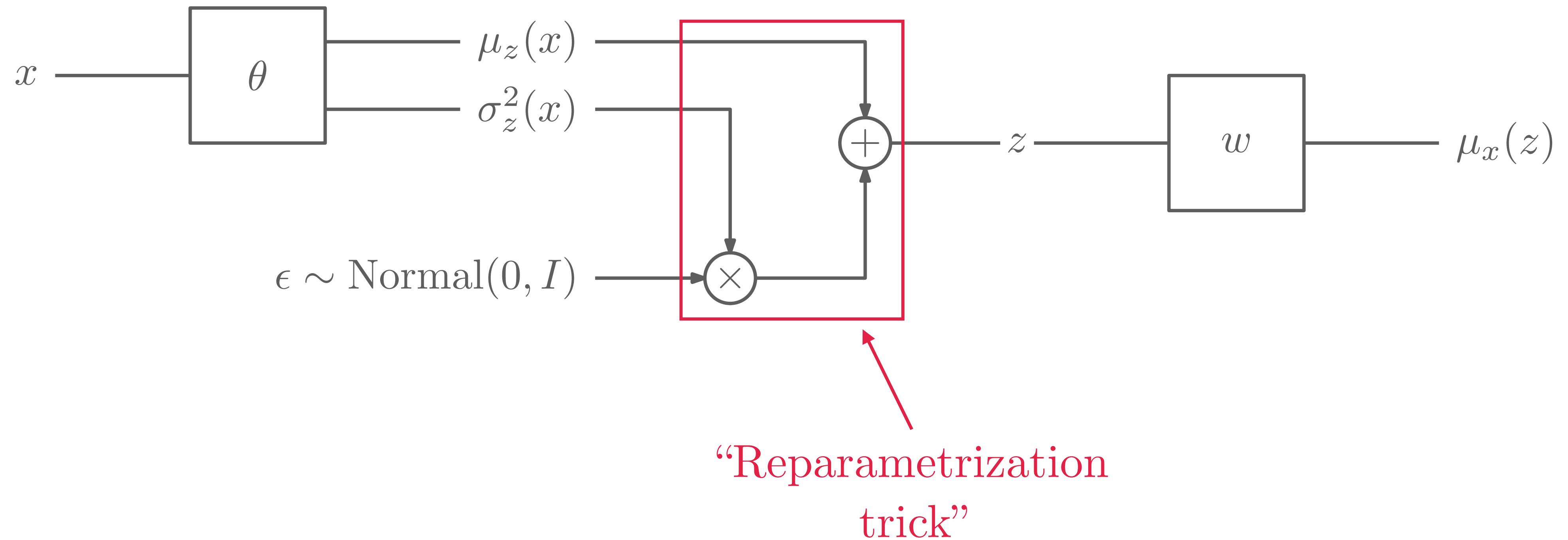
with $z \sim q_\theta(\cdot | x)$

Training a VAE



How do we backprop
through sampling??

Training a VAE



Reconstruction using VAE

7|0|3|1|2|7|0|2|9|6|0|3|1|6|7|1|9|7|6|5|5|4|8|3|4|4|8|7|3|6

Original
images

7|0|3|1|2|7|0|2|9|6|0|3|1|6|7|1|9|7|6|5|5|4|8|3|4|4|8|7|3|6

Reconstructed
images

Reconstruction using VAE



Original
images

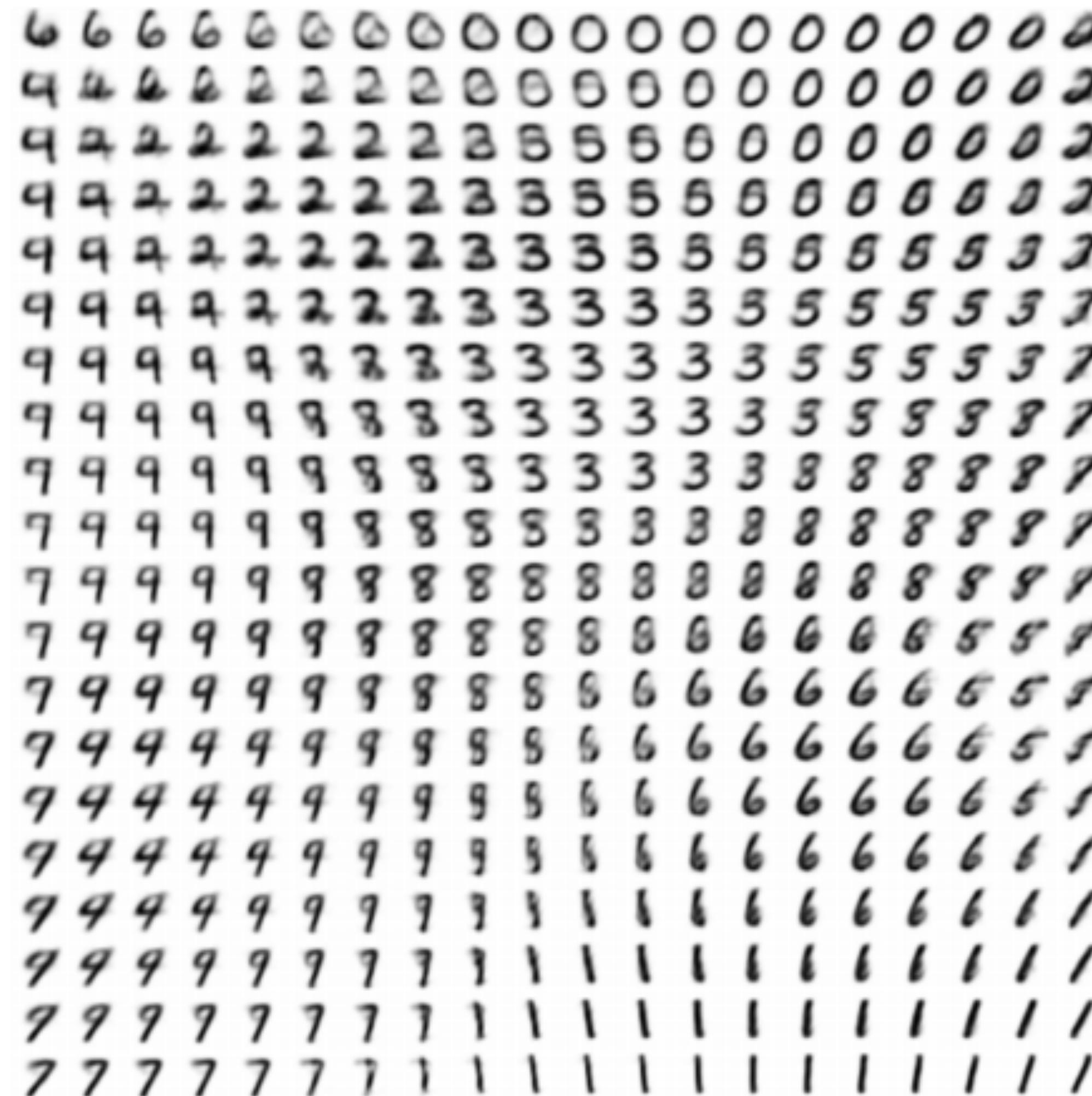


Reconstructed
images

Generation using a VAE

- We can now use the trained decoder to generate artificial samples:
 - We sample $z \sim p_{\text{prior}}$
 - We use the decoder to generate synthetic samples \hat{x} from z , $\hat{x} \sim p_w(\cdot | z)$

Generation using a VAE



Generated images
as latent space is
traversed

Generation using a VAE



Generated images
as latent space is
traversed

How can we train p_w ?
(reprise)

Training p_w

- VAEs depart from a Bayesian approach to training p_w
 - The approach requires “friendly” distributions (e.g., Gaussians)
 - The resulting images are “blurry”
 - VAEs may suffer from **posterior collapse** (i.e., learning to ignore latent variable)

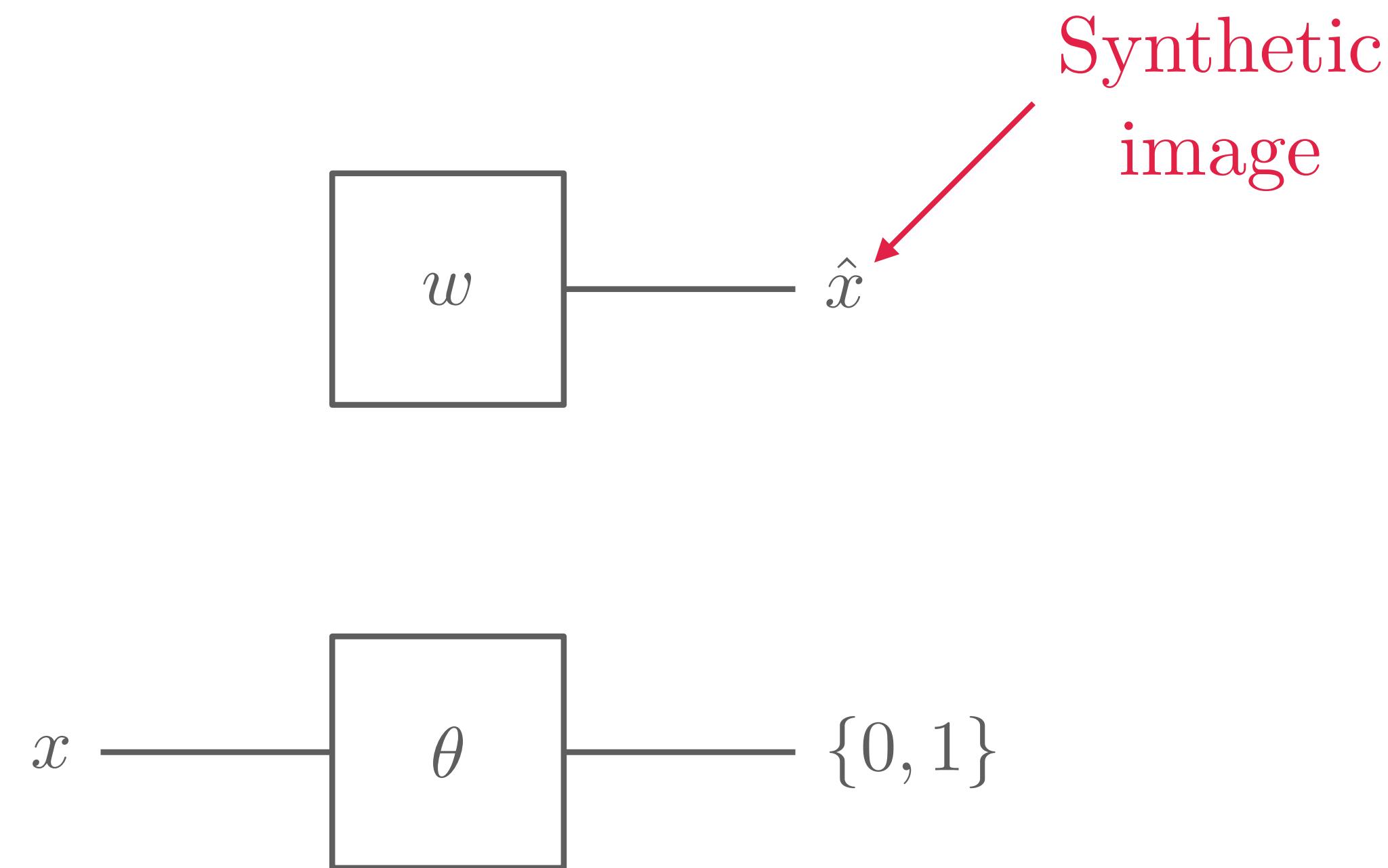
Training p_w

- If the goal is to generate realistic data, why not evaluate the data on “realism”?
 - Idea:
 - Train our generator network to produce realistic synthetic data
 - Train a second network (discriminator) to distinguish real from synthetic data
- Performance: Ability to fool
the discriminator
- Performance: Ability to
recognize real from
fake images

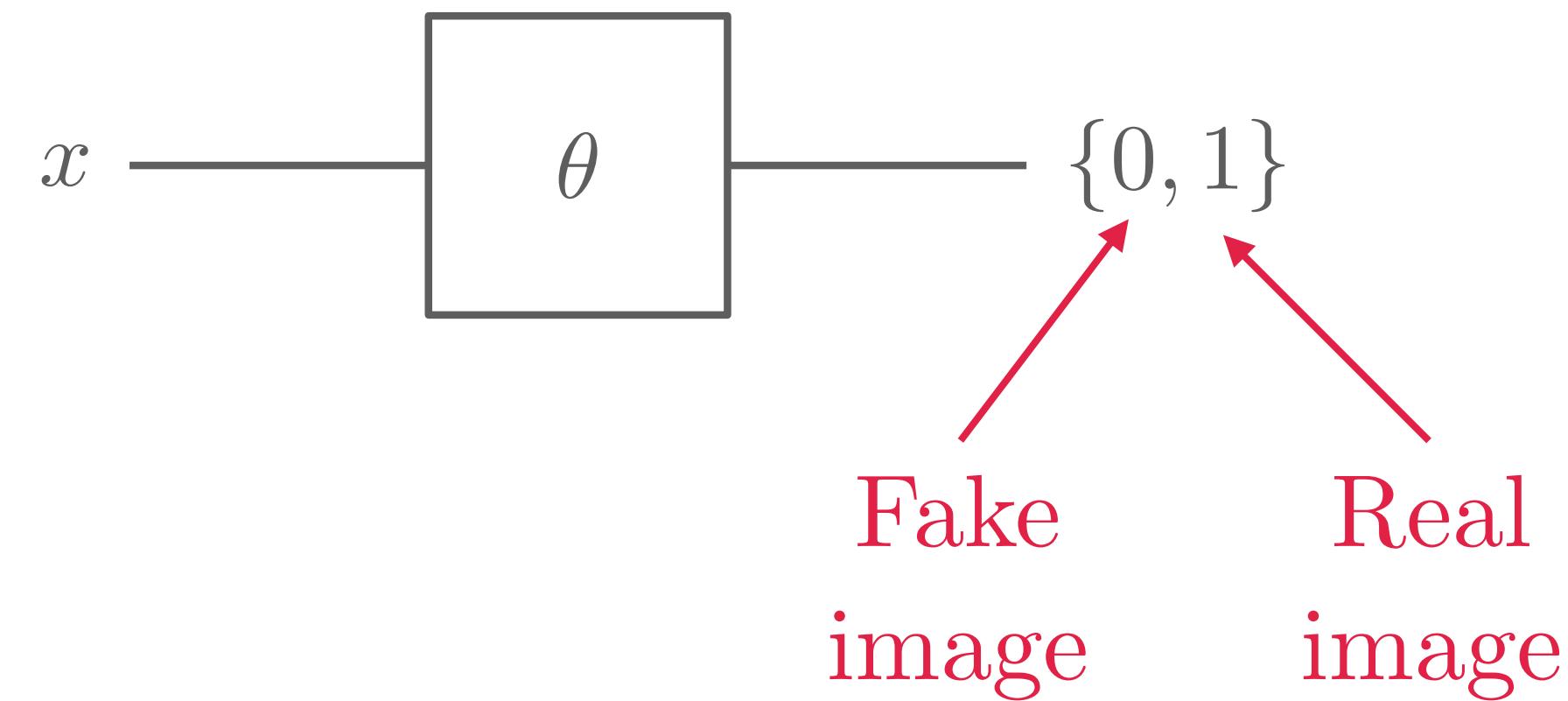
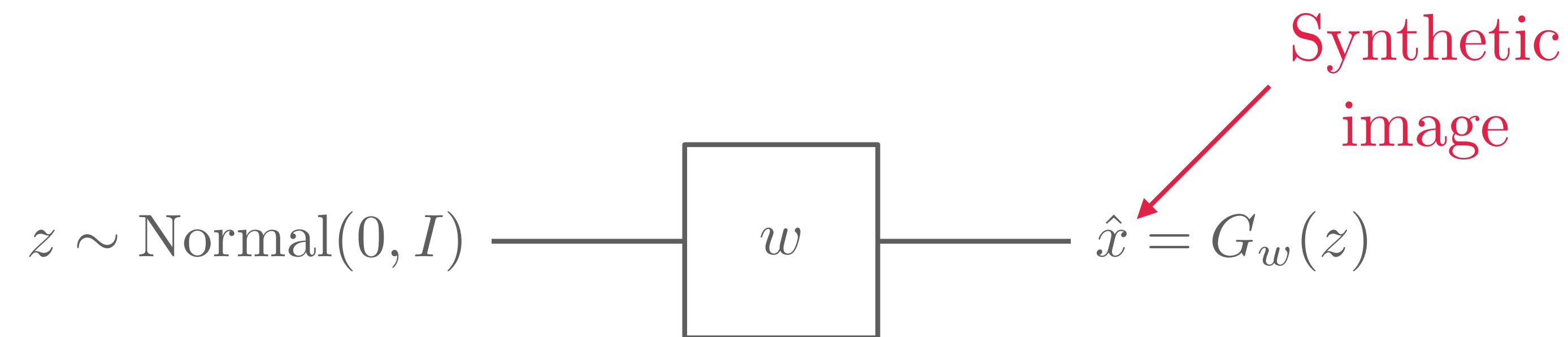
Training p_w

- If the goal is to generate realistic data, why not evaluate the data on “realism”?
- Idea:
 - Train our generator network to produce realistic synthetic data
 - Train a second network (discriminator) to distinguish real from synthetic data
 - The training process is **adversarial** (the two networks are “competing” with one another)

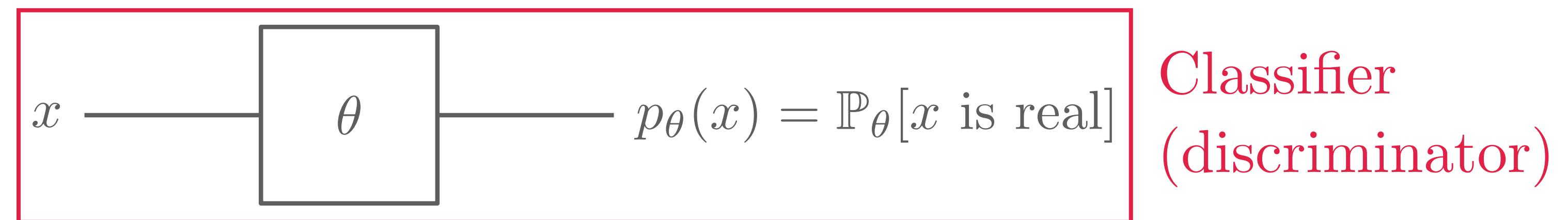
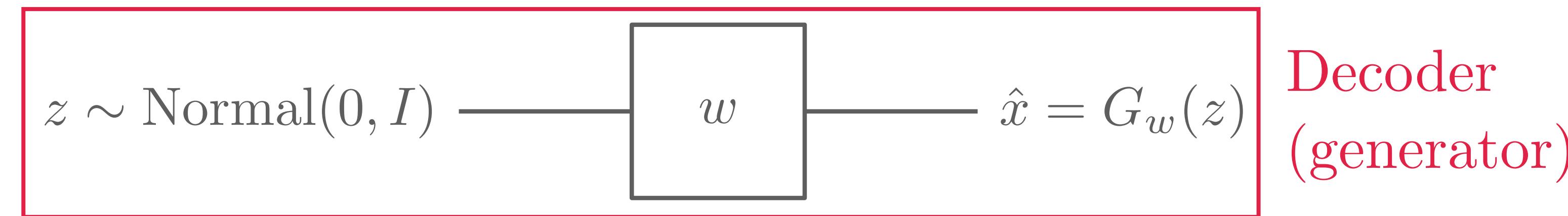
Generative adversarial network (GAN)



Generative adversarial network (GAN)



Generative adversarial network (GAN)



Generative adversarial network (GAN)

- Generator network, G , and discriminator network, D , are playing a **game**
 - Discriminator wants to **maximize**

$$\mathbb{E}_{x \sim \mu_{\text{data}}} [\log p_{\theta}(x)]$$

Real
data

$$\mathbb{E}_{x \sim p_w(\cdot|z)} [\log p_{\theta}(x)]$$

“Fake”
data

Generative adversarial network (GAN)

- Generator network, G , and discriminator network, D , are playing a **game**
 - Discriminator wants to **maximize**

$$\mathbb{E}_{x \sim \mu_{\text{data}}} [\log p_{\theta}(x)]$$

and **minimize**

$$\mathbb{E}_{z \sim \text{Normal}(0, I)} [\log p_{\theta}(G_w(z))]$$

Generative adversarial network (GAN)

- Generator network, G , and discriminator network, D , are playing a **game**
 - Discriminator wants to **maximize**

$$\mathbb{E}_{x \sim \mu_{\text{data}}} [\log p_{\theta}(x)]$$

and

$$\mathbb{E}_{z \sim \text{Normal}(0, I)} [\log(1 - p_{\theta}(G_w(z)))]$$

Generative adversarial network (GAN)

- Generator network, G , and discriminator network, D , are playing a **game**
 - Generator wants to **minimize**

$$\mathbb{E}_{x \sim \mu_{\text{data}}} [\log p_{\theta}(x)]$$

and

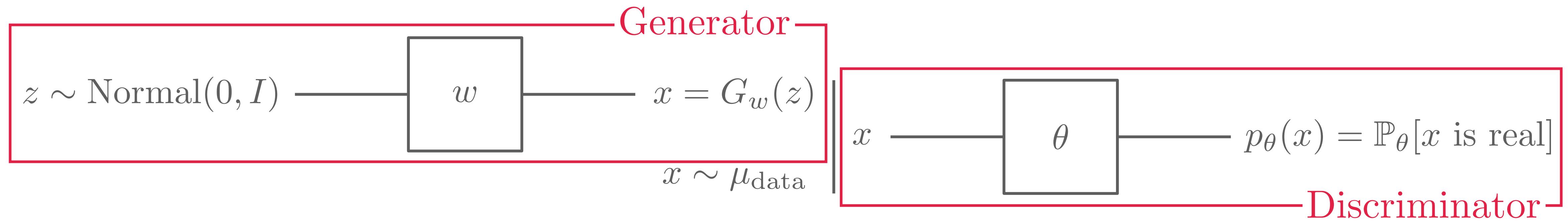
$$\mathbb{E}_{z \sim \text{Normal}(0, I)} [\log(1 - p_{\theta}(G_w(z)))]$$

Generative adversarial network (GAN)

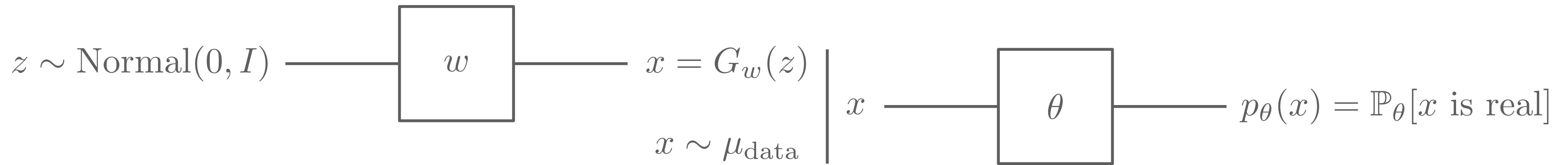
- Generator network, G , and discriminator network, D , are playing a **game**
- Saddle point problem:

$$\min_w \max_{\theta} \mathbb{E}_{x \sim \mu_{\text{data}}} [\log p_{\theta}(x)] + \mathbb{E}_{z \sim \text{Normal}(0, I)} [\log(1 - p_{\theta}(G_w(z)))]$$

Putting everything together...



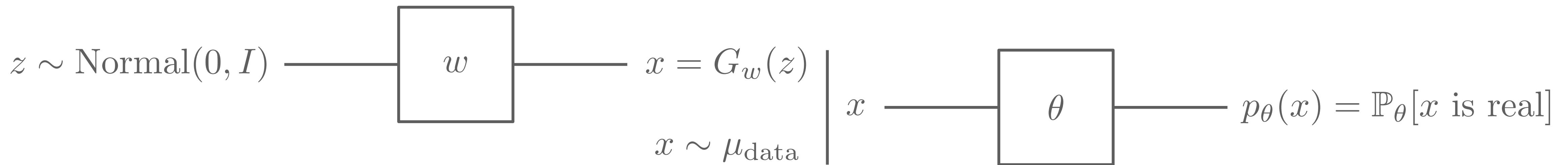
Training a GAN



- We want to solve the minimax problem

$$\min_w \max_\theta \mathbb{E}_{x \sim \mu_{\text{data}}} [\log p_\theta(x)] + \mathbb{E}_{z \sim \text{Normal}(0, I)} [\log(1 - p_\theta(G_w(z)))]$$

Training a GAN



- We use stochastic gradient descent
- We alternate between:

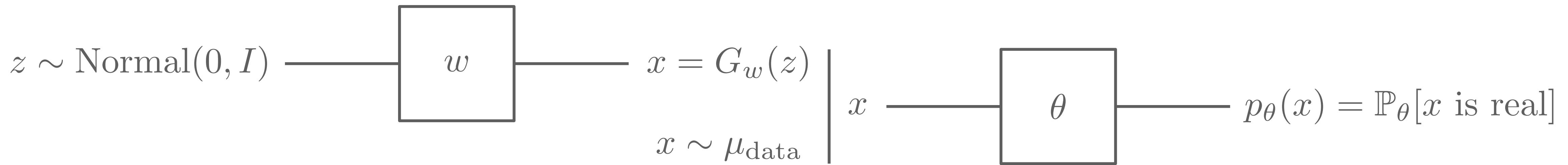
Minibatch

- Minimizing the generator loss,

$$L_G(w) = \frac{1}{N} \sum_{n=1}^N \log(1 - p_\theta(G_w(z_n)))$$

size
 $z_n \sim \text{Normal}(0, I)$

Training a GAN



- We use stochastic gradient descent
- We alternate between:
 - ... and minimizing the **discriminator loss**,

$$L_D(\theta) = -\frac{1}{N} \sum_{n=1}^N \log p_\theta(x_n) - \log(1 - p_\theta(G_w(z_n)))$$

Sample from
dataset

$$z_n \sim \text{Normal}(0, I)$$

Training a GAN

- Several variants and schedules exist
- In general,
 - The optimization of GANs is often difficult
 - There are no convergence guarantees

Generation using GANs



8916393150
2909804320
9712473064
3911813496
1705457042
5791749812
7011320143
0854809681
8855549966
2090061999

Generation using GANs



Generation using GANs



Concluding remarks

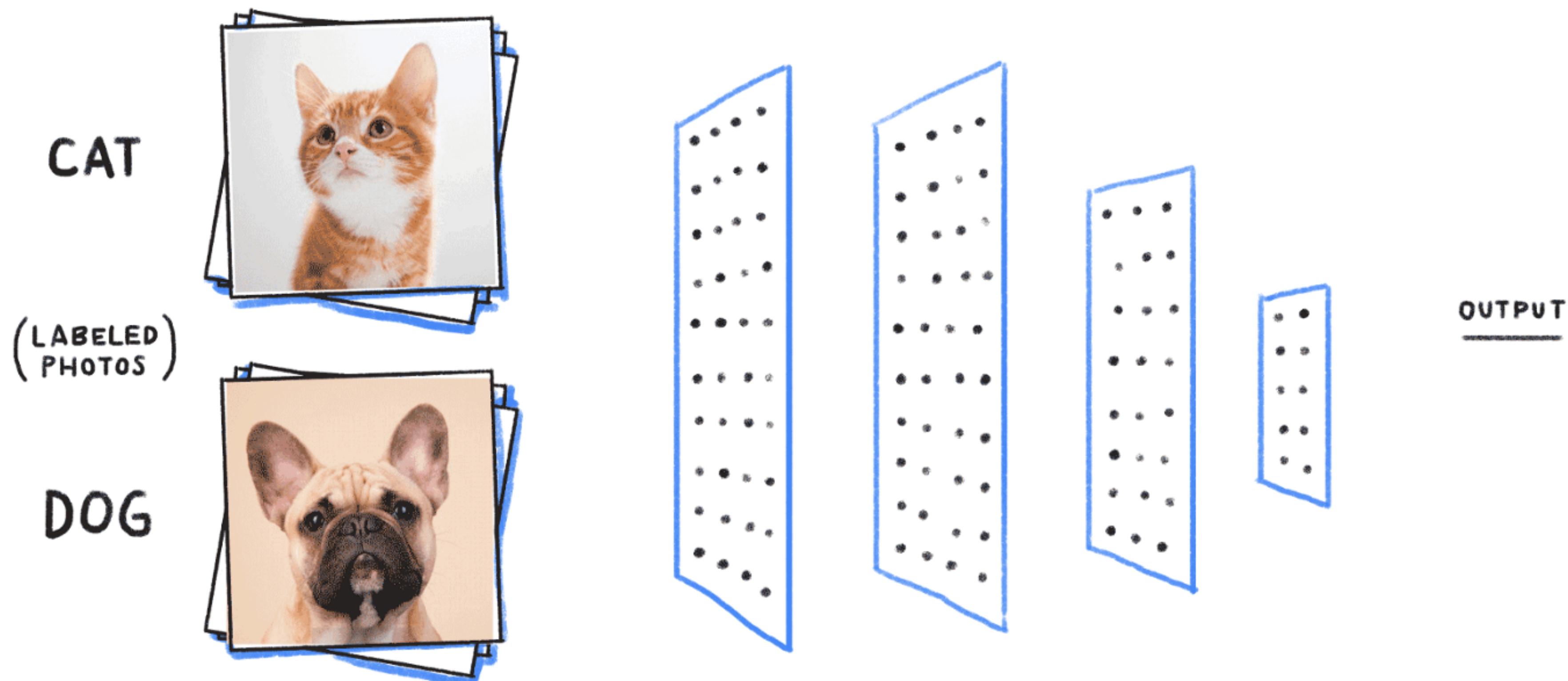
- GANs have several points in their favor:
 - GANs generate the sharpest images
 - They are cheaper to train than other generative models
- However, ...
 - They are difficult to optimize (often unstable training dynamics)
 - We cannot use them to perform statistical inference

The ~~other~~ side of deep learning

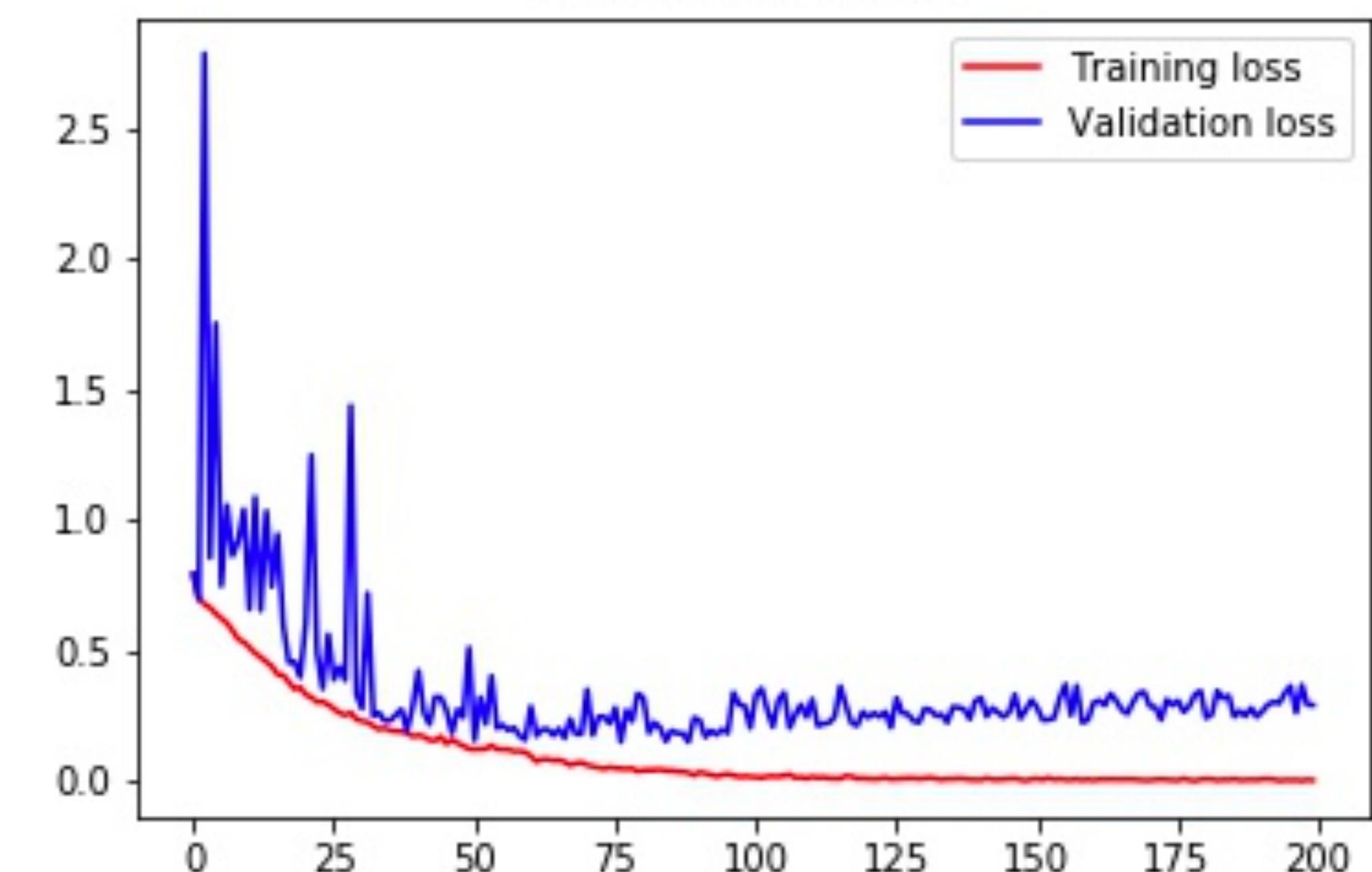
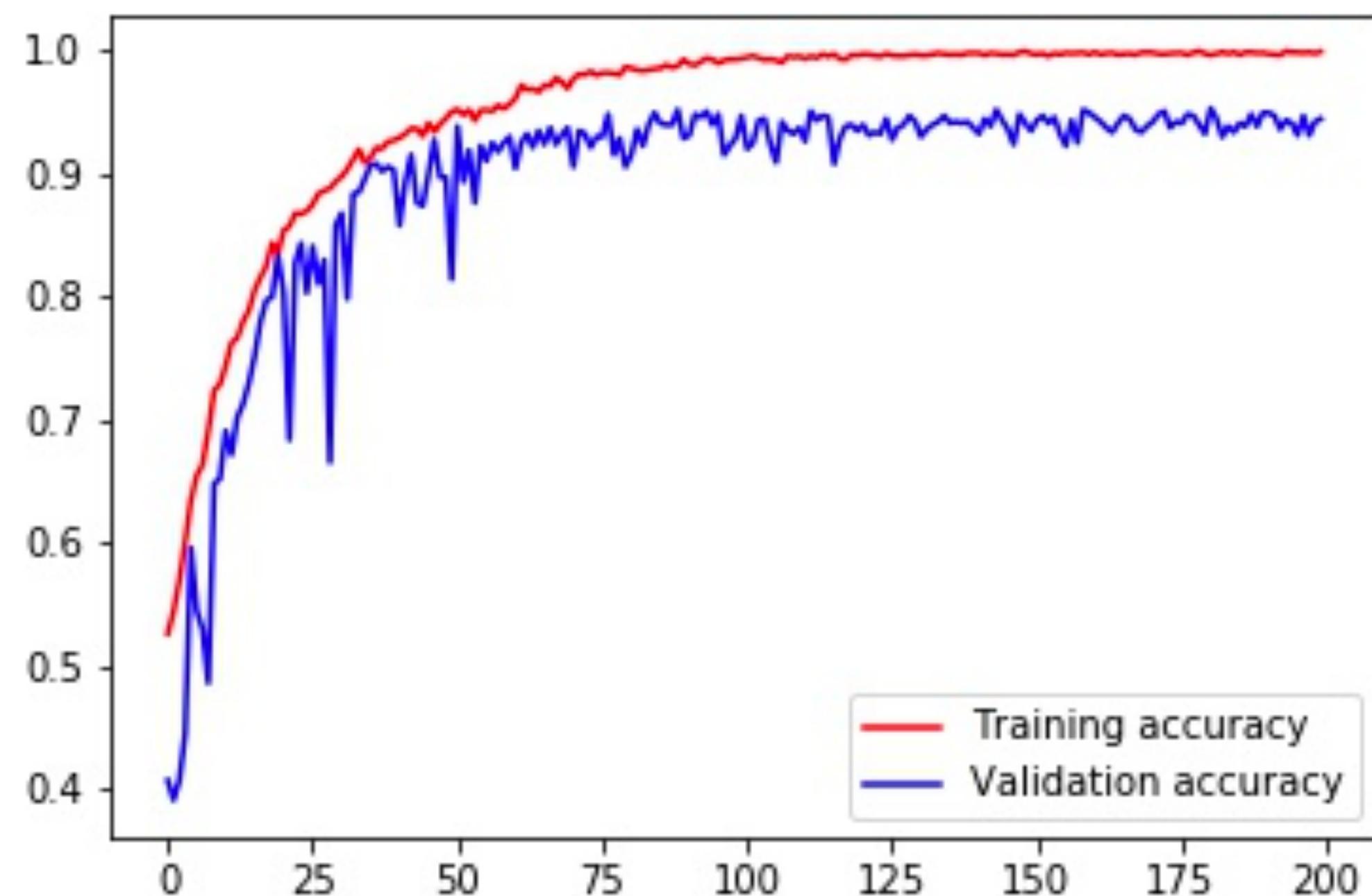
What other side?

- Are deep learning models trustworthy?
- Can we explain the outputs of a deep learning model?
- Fairness and ethics

Example

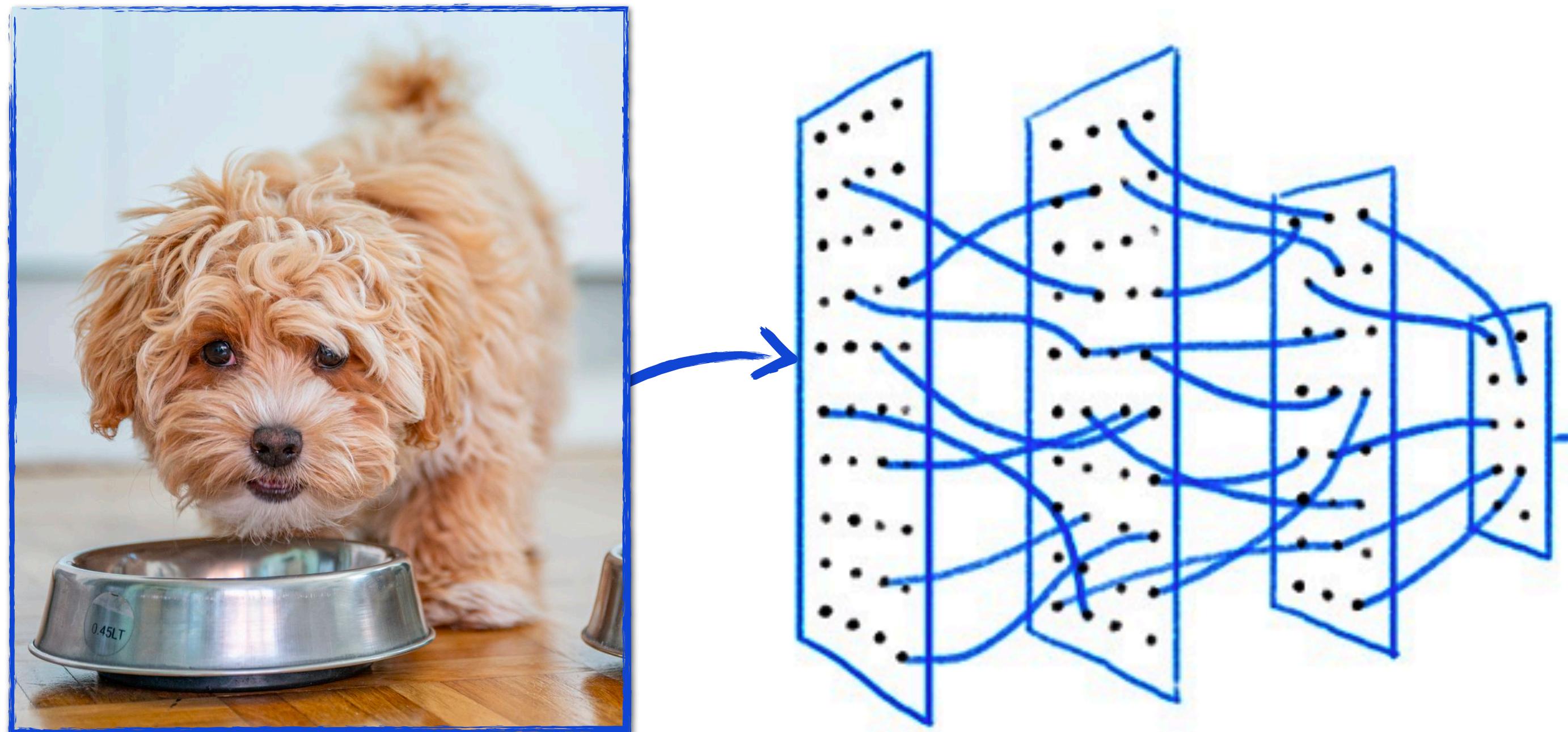


Example

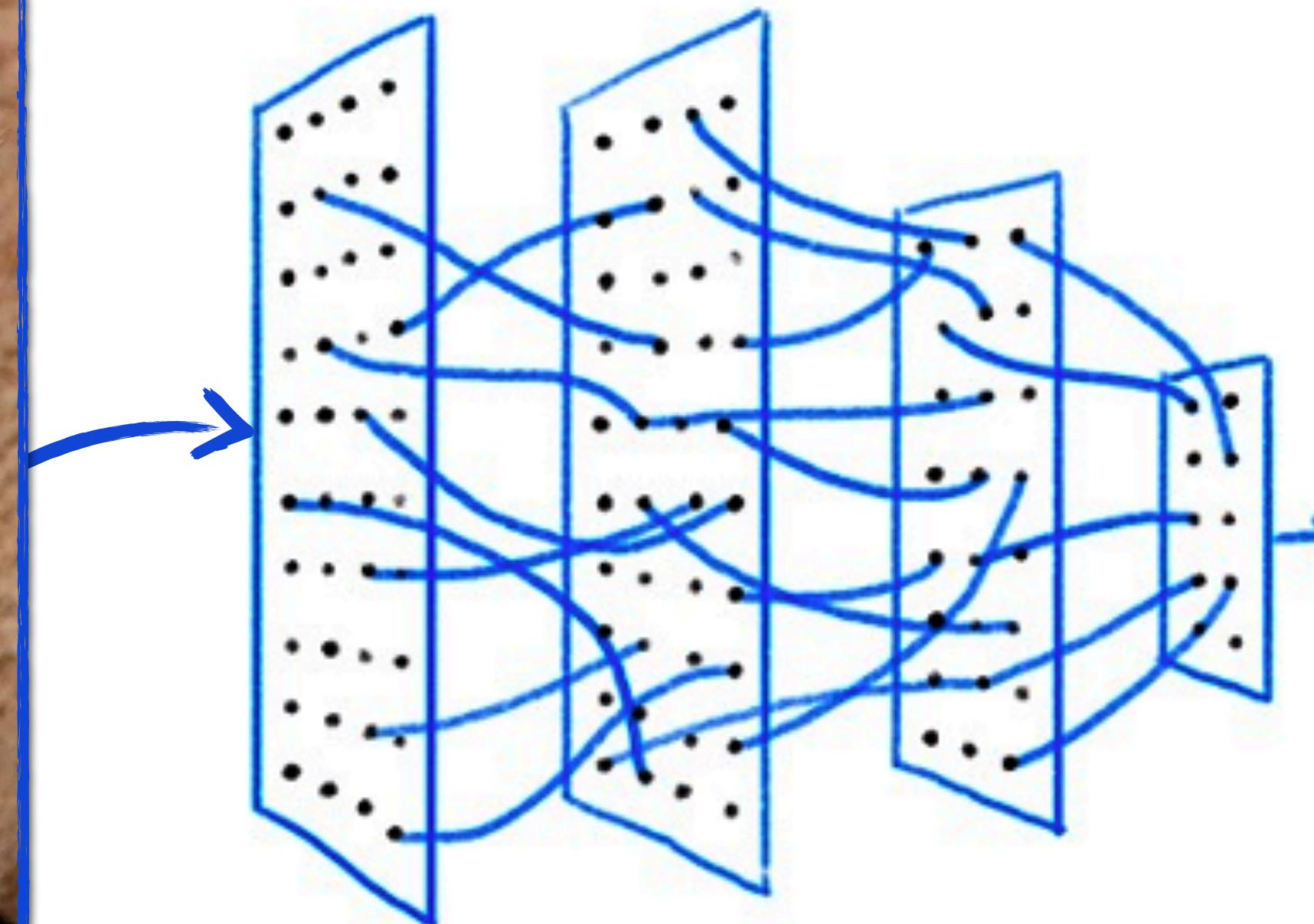
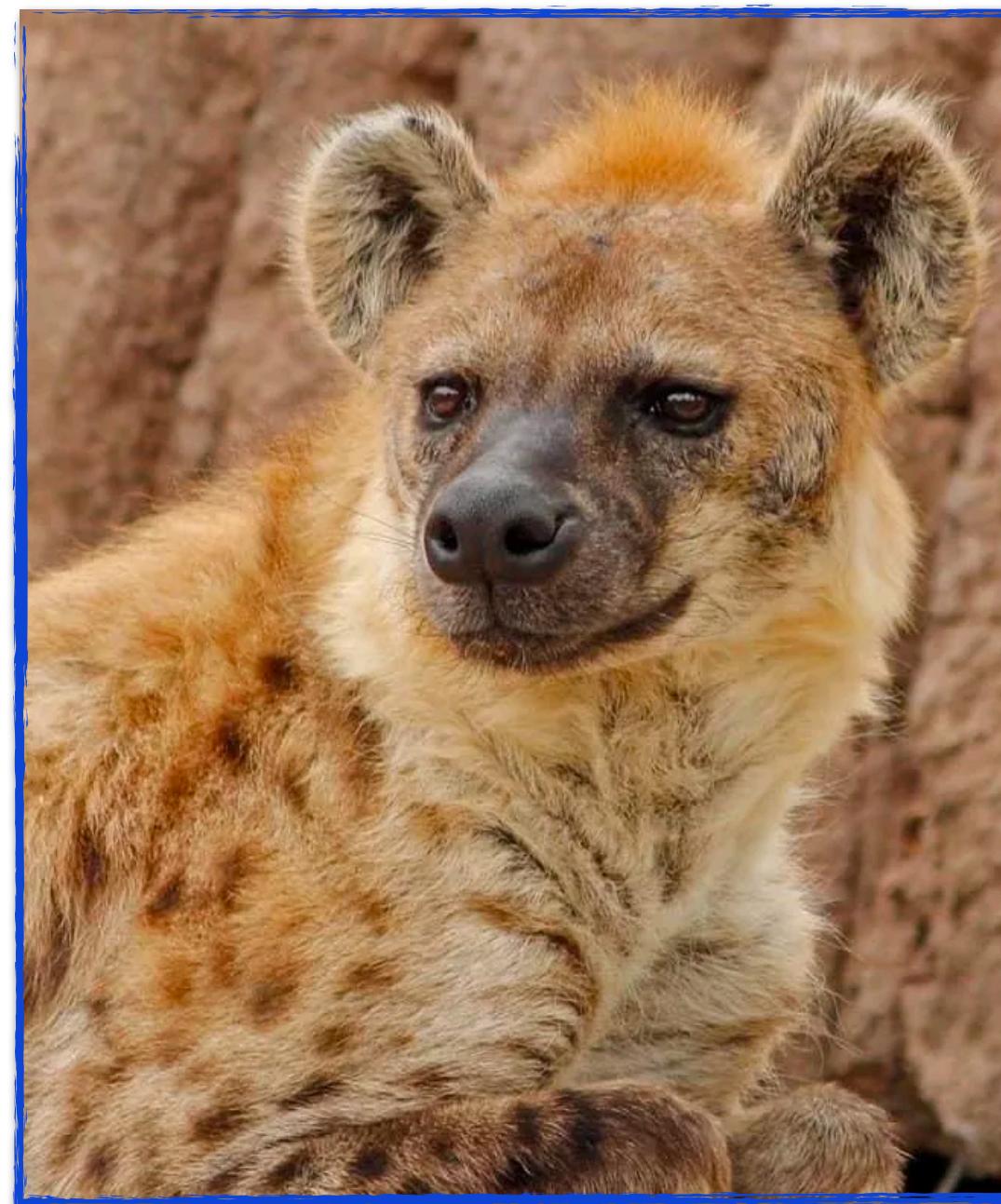


Would you **trust** this model?

Can we trust deep learning models?

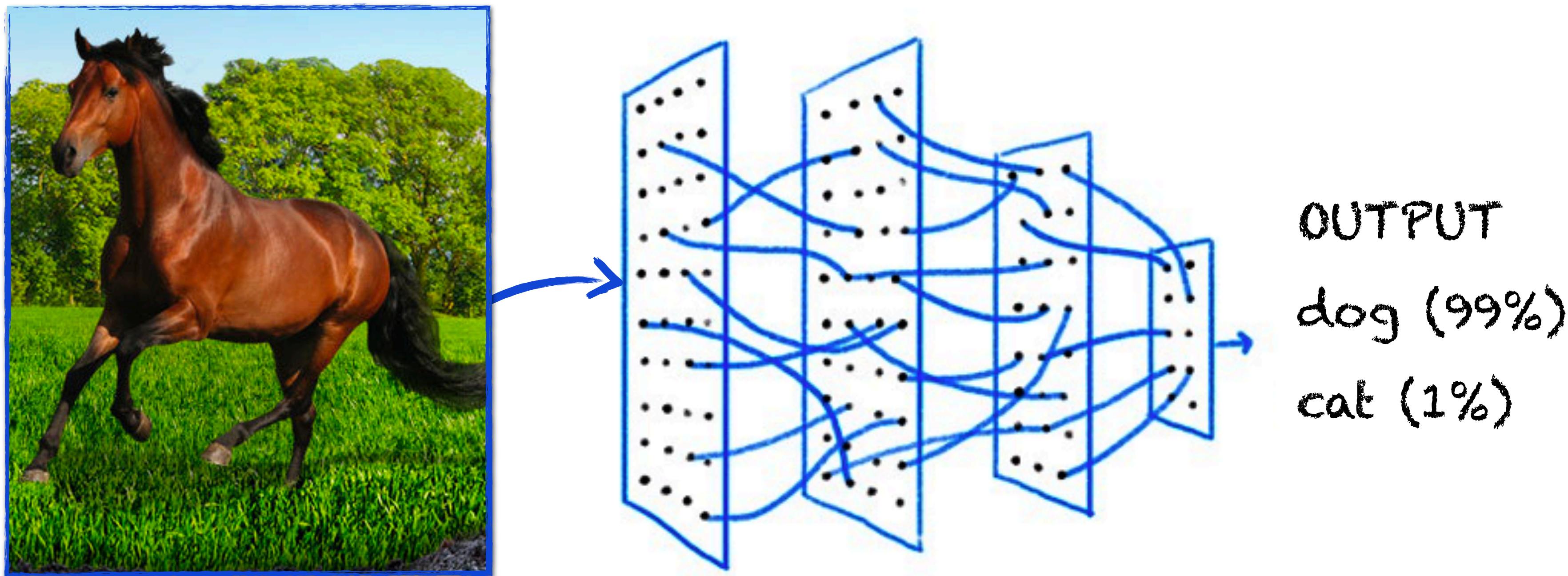


Can we trust deep learning models?



OUTPUT
dog (90%)
cat (10%)

Can we trust deep learning models?



Can we trust deep learning models?

- The model may be expected to be wrong in **out-of-domain** inputs
- Can we have some indication about when the model is “confused”?



Uncertainty

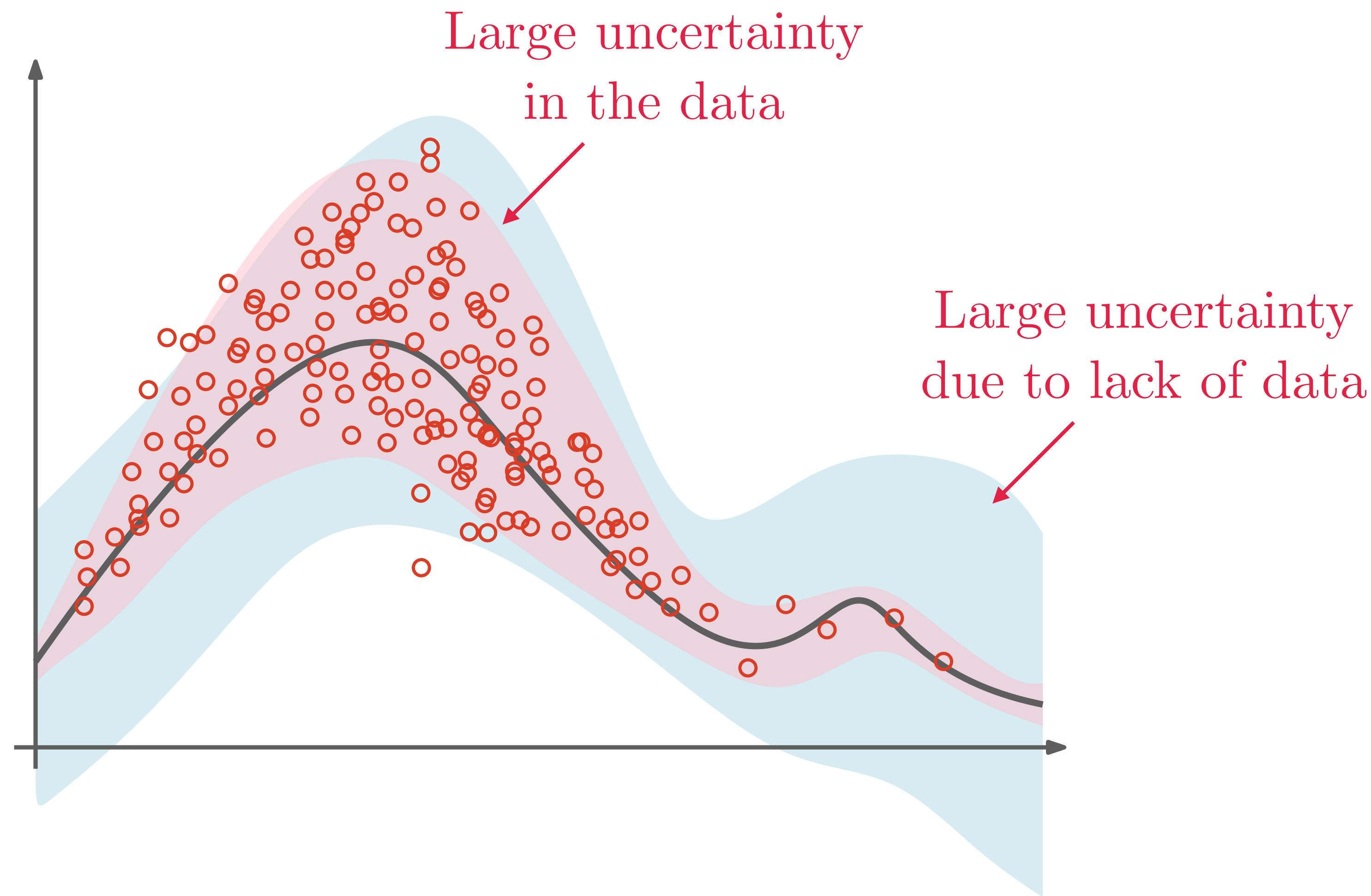
Uncertainty in DL predictions

- Assessing uncertainty in DL predictions is useful for...
 - Deciding when a prediction must be further inspected
 - Debugging, retraining, etc.
 - Providing safe-checks for high-risk applications

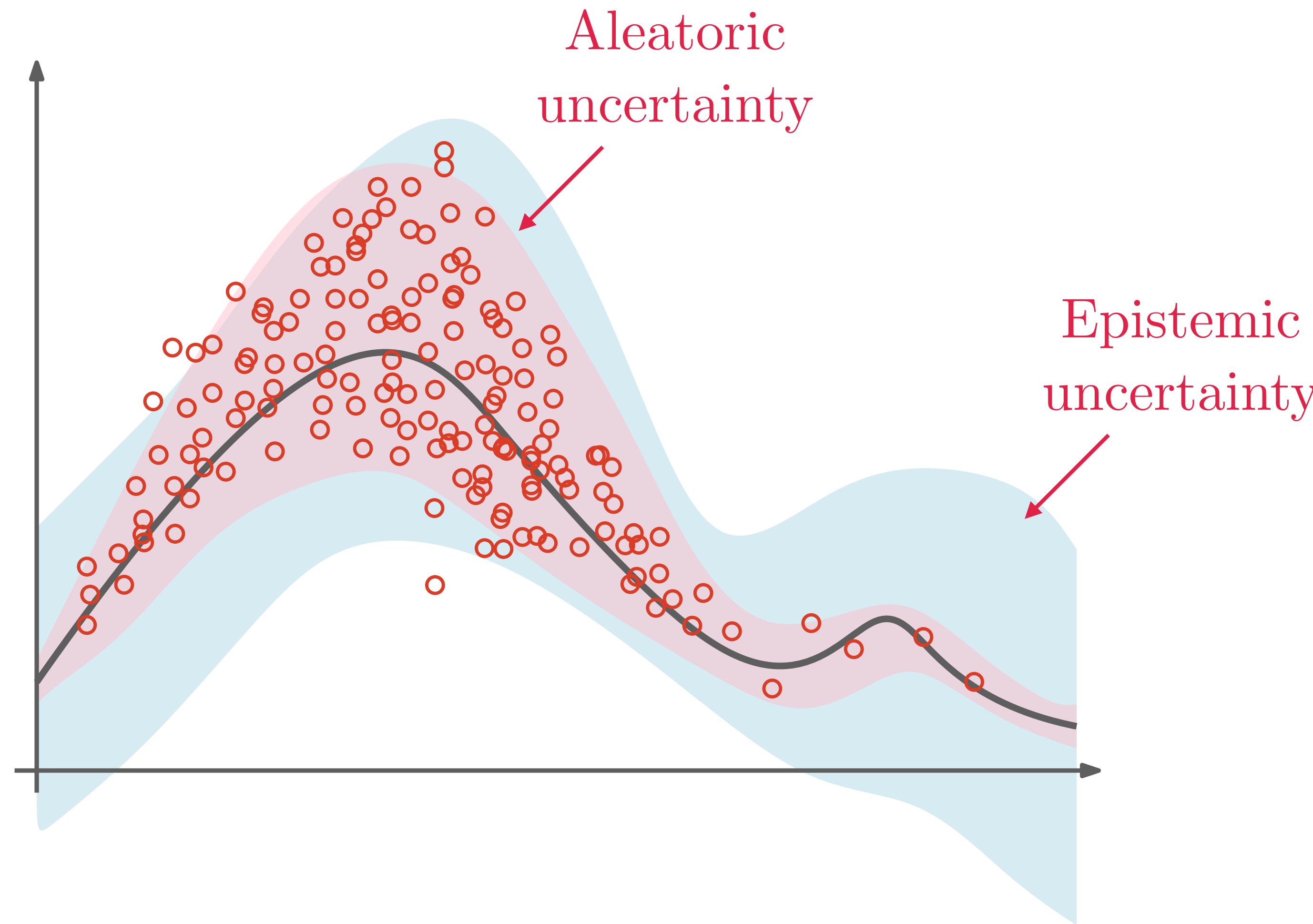
Uncertainty in DL predictions

- Uncertainty in the output of a deep learning model may be caused by:
 - ... uncertainty inherent to the data (**aleatoric uncertainty**)
 - ... knowledge/model uncertainty (**epistemic uncertainty**)

Uncertainty in DL predictions



Uncertainty in DL predictions



Uncertainty in DL predictions

Measure of
“confusion”

- Estimating aleatoric uncertainty:
 - We can estimate data **entropy** in the dataset
 - We can use data from **multiple annotators** to estimate uncertainty

Uncertainty in DL predictions

- Estimating epistemic uncertainty:
 - Often due to lack of data / out-of-domain data
 - How much do predictions change with data?



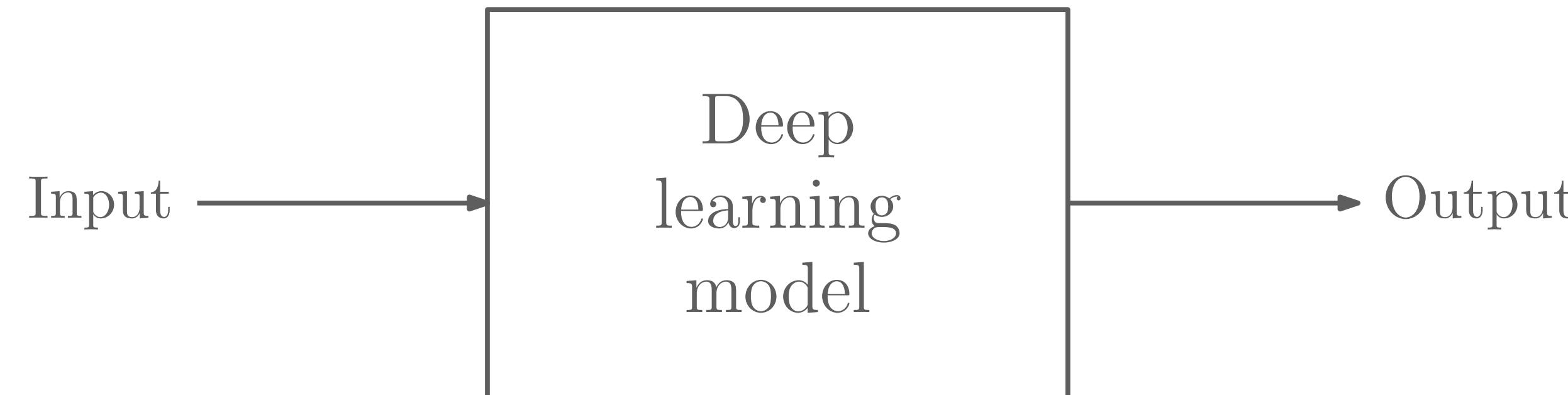
Variance

Variance estimation in DL models

- Estimating variance:
 - Monte-Carlo Dropout (use dropout at estimation time)
 - Deep ensembles (independently training multiple models)
- ...

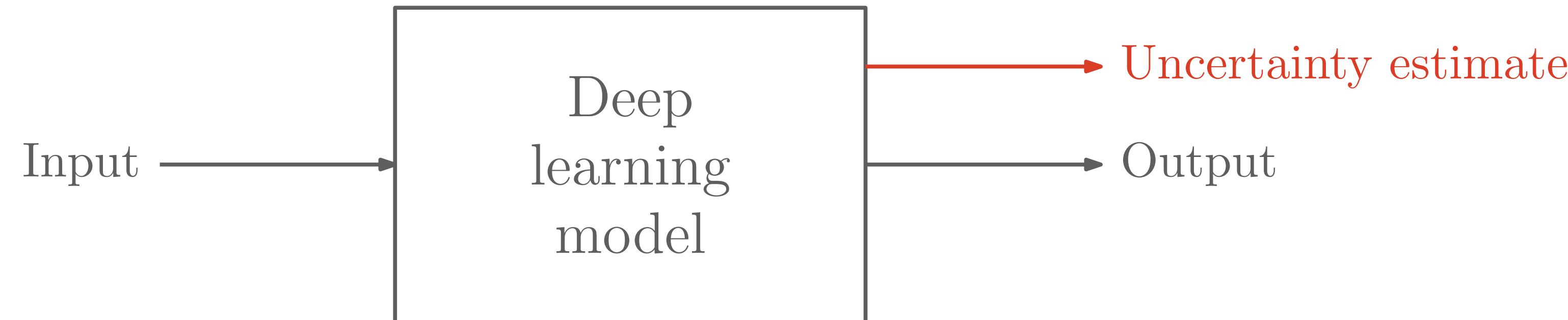
Beyond uncertainty

- Uncertainty provides a way of quantifying to what extent predictions can be trusted



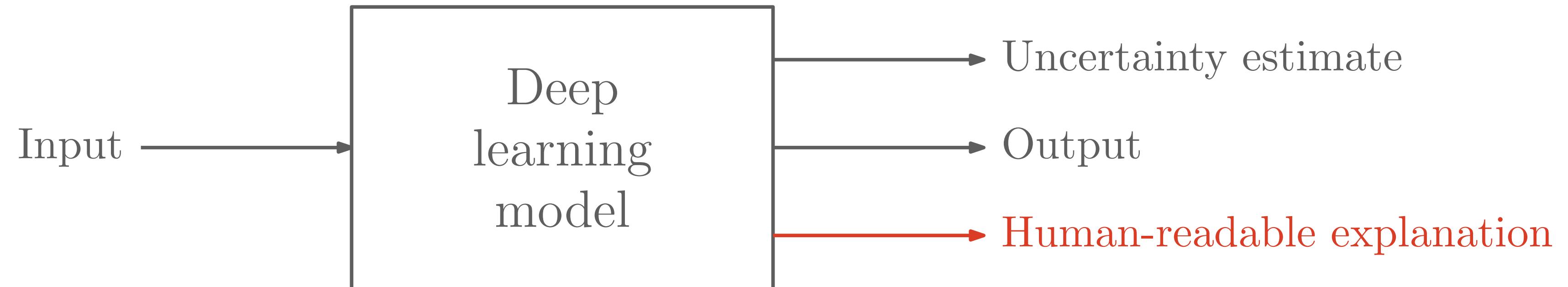
Beyond uncertainty

- Uncertainty provides a way of quantifying to what extent predictions can be trusted



Beyond uncertainty

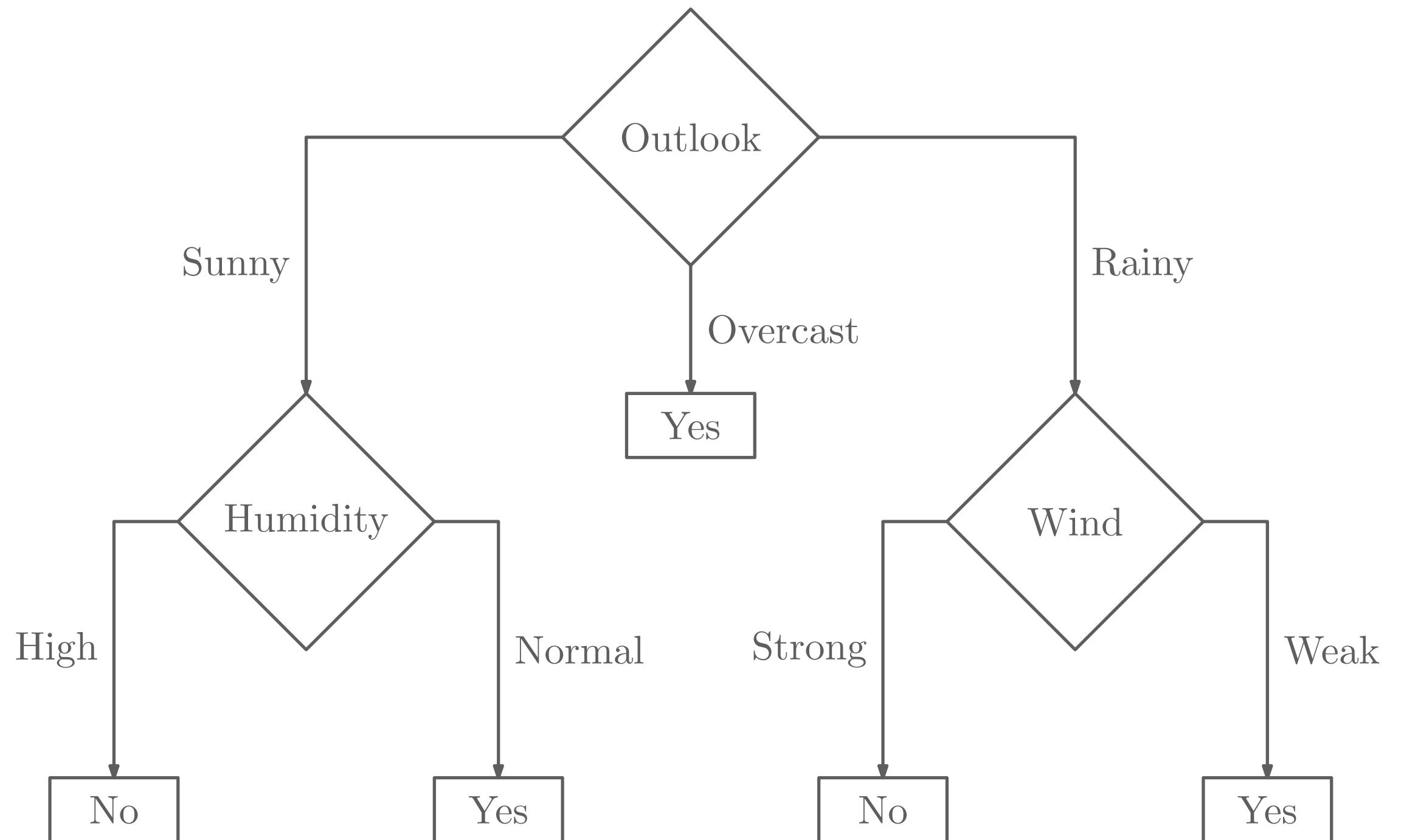
- Uncertainty provides a way of quantifying to what extent predictions can be trusted
- However, it is often not enough to understand (and trust) the predictions made by DL models



Explainability in machine learning

- Some machine learning models are **intrinsically explainable**
- E.g.,
 - Decision trees

Can turn it into a set of “if-then” rules



Explainability in machine learning

- Some machine learning models are **intrinsically explainable**
- E.g.,
 - Decision trees

If the outlook is **overcast**, then **yes**.

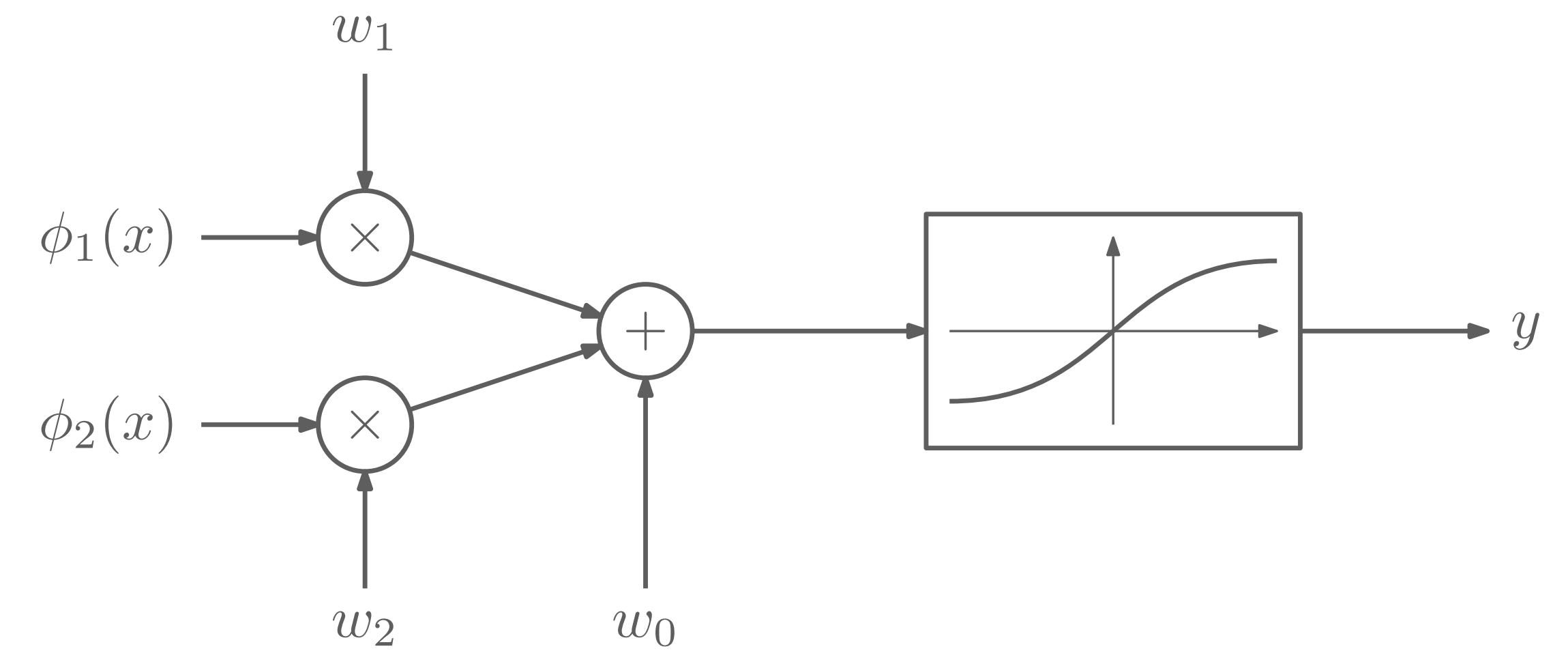
Similarly, if it's **sunny** but the **humidity** is **normal**, then **yes**.

Same thing if it's **rainy** but the **wind** is **weak**, then **yes**.

Otherwise, **no**.

Explainability in machine learning

- Some machine learning models are **intrinsically explainable**
- E.g.,
 - Decision trees
 - Linear models



$$y = F(w_0 + w_1\phi_1(x) + w_2\phi_2(x))$$

Relative importance
of different features

Can we have post-hoc explanations?

- We can “build” local explanations for a trained model

Explanations computed
for specific inputs



Can we have post-hoc explanations?

- We can “build” local explanations for a trained model
 - Model-specific explanations (can be applied to specific models)
 - Model-agnostic explanations (can be applied to “arbitrary” models)

Model-specific explanations

- Gradient-based explanations
 - Can be used in **differentiable** machine learning models (e.g., neural networks)
 - E.g., use gradients to reveal which regions of input contribute the most to the prediction

$$F_w(x + \epsilon) \approx F(x) + \epsilon^\top \nabla_x F_w(x) + \text{small stuff}$$


Gradient w.r.t. input indicates
which input directions (features)
most impact the output

Model-specific explanations

- Attention-based explanations
 - Can be used with attention-based models
 - E.g., use attention weights to determine most relevant inputs

fantastic movie one of the best film noir movies ever made

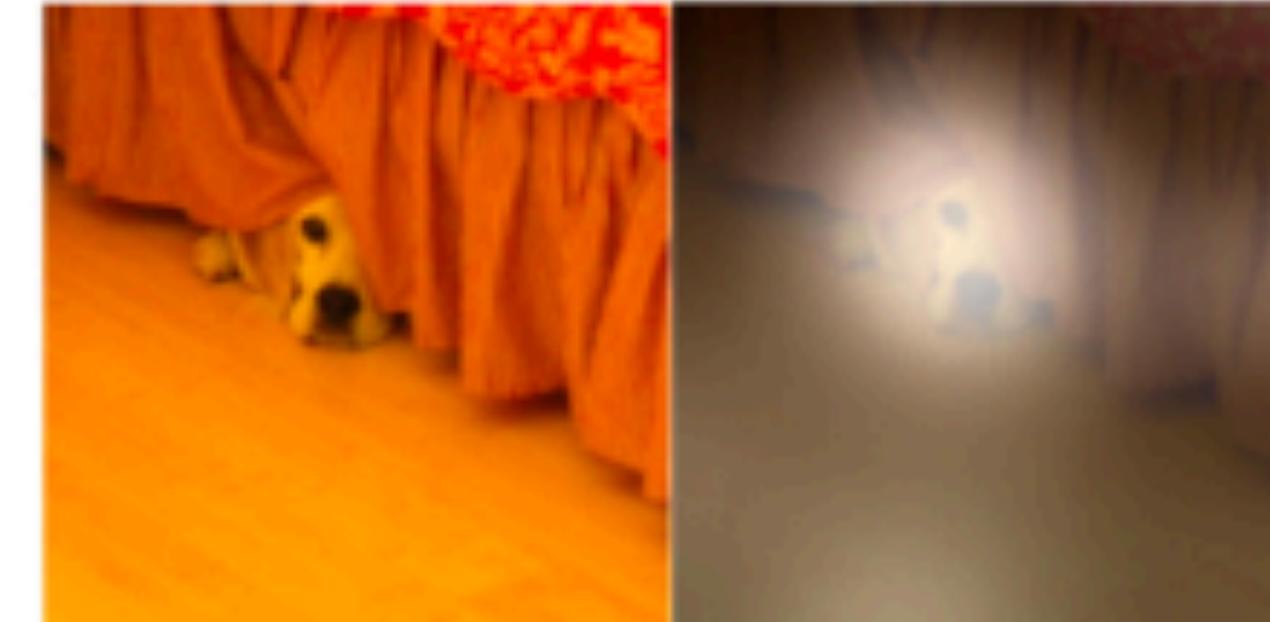
a meandering inarticulate and ultimately disappointing film

Model-specific explanations

- Attention-based explanations
 - Can be used with **attention-based** models
 - E.g., use **attention weights** to determine most relevant inputs



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.

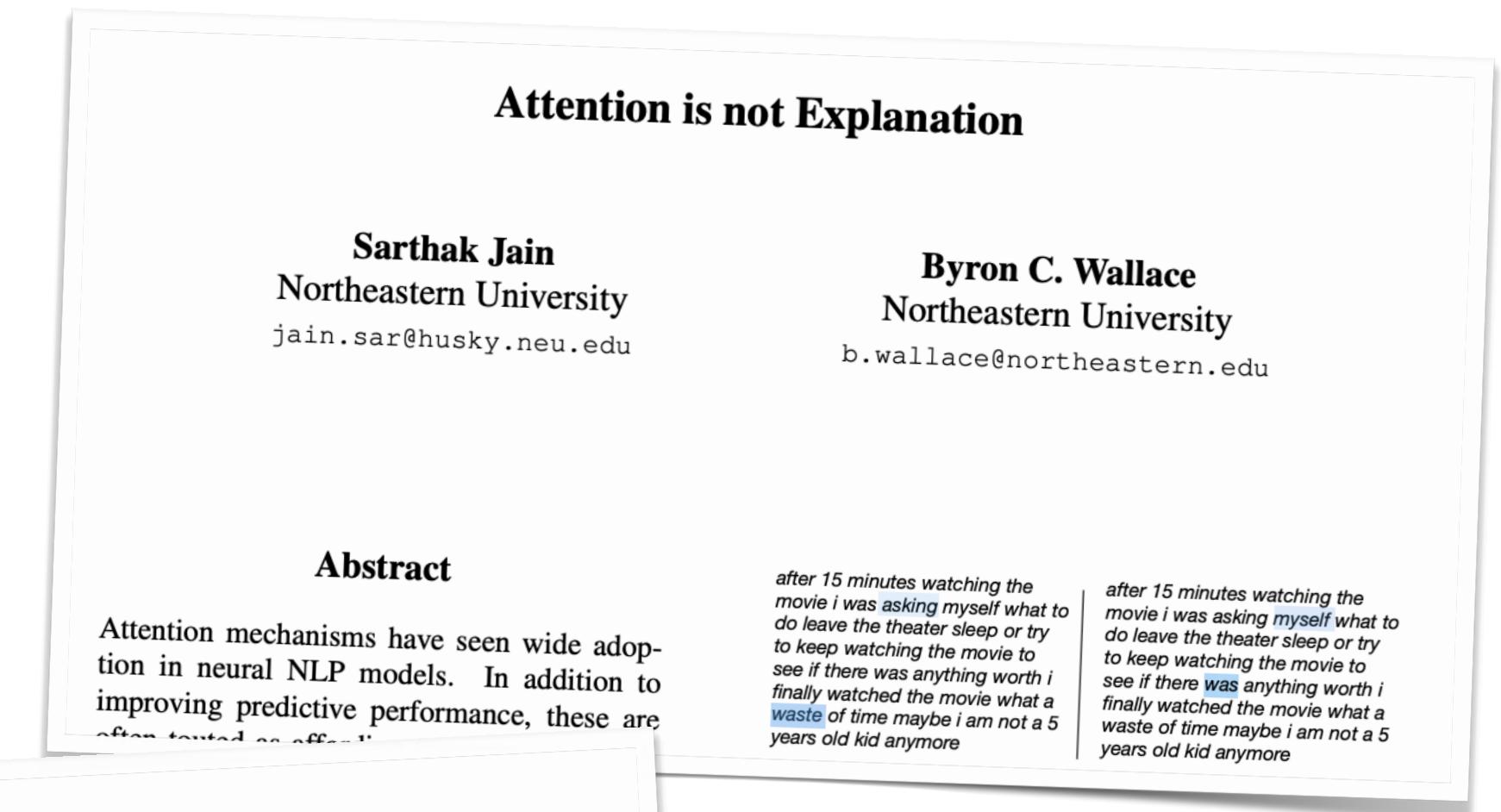
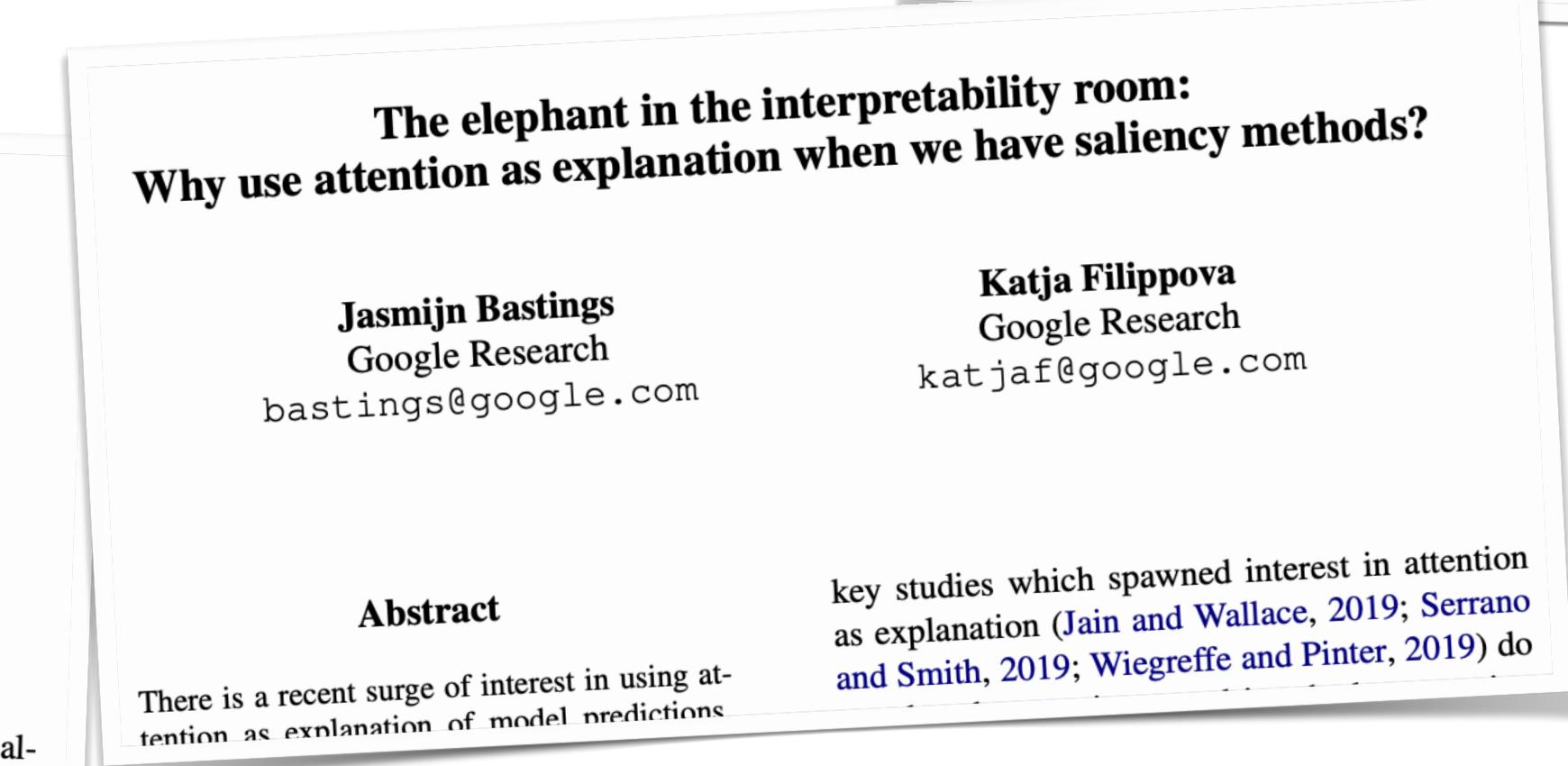
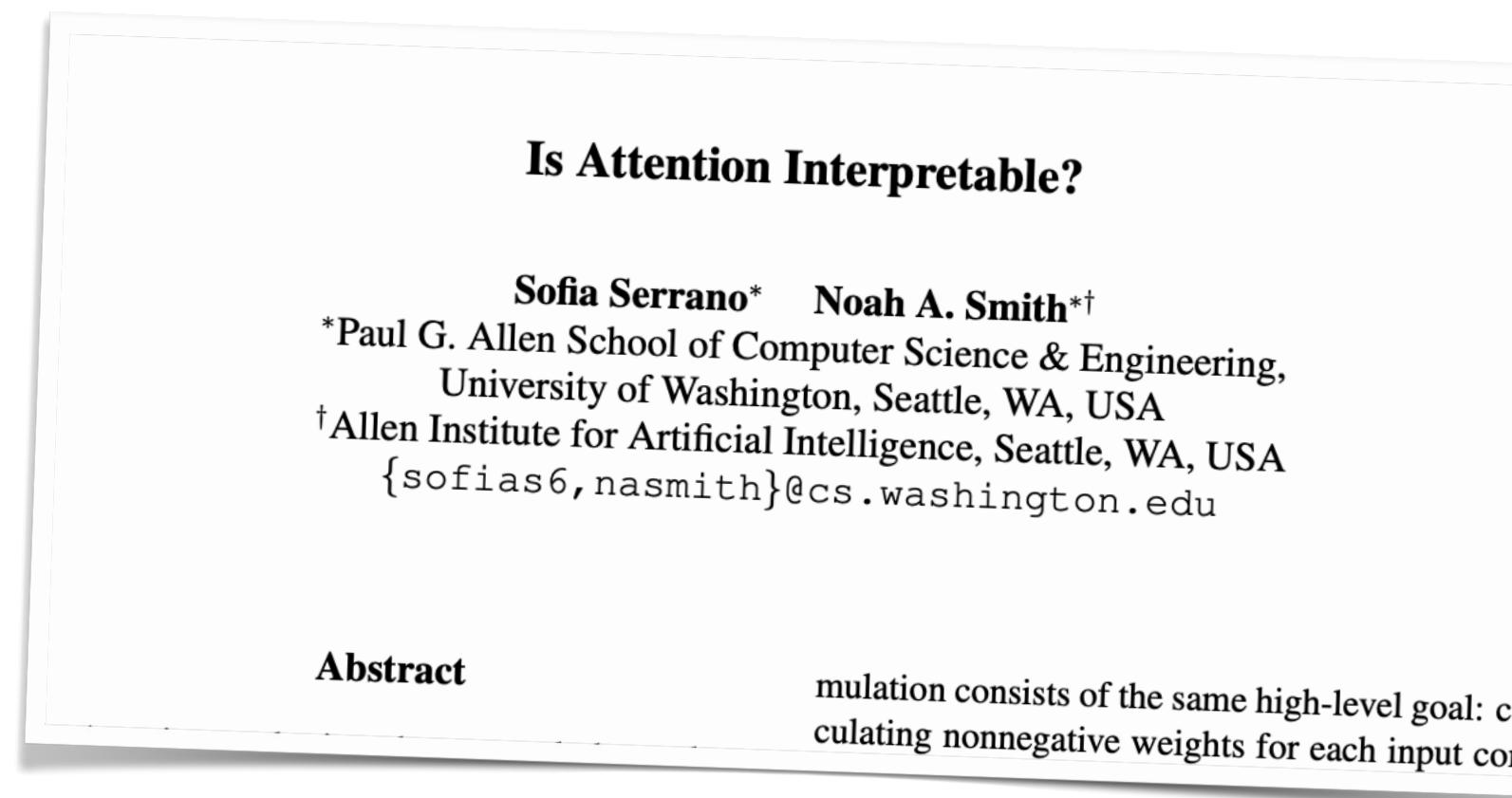


A stop sign is on a road with a mountain in the background.

Model-specific explanations

- Attention-based explanations

- Can be used with attention-based models
- Not universally accepted...



Model-agnostic explanations

- Perturbation-based explanations
 - Share intuition with gradient-based explanations
 - Seek to answer the question:
“Which part of the input most influences the prediction?”
 - E.g., LIME (Local Intepretable Model-agnostic Explanations)

LIME

- For given input x ,
 - Select features of the input

Words of
a sentence

Regions of
an image

LIME

- For given input x ,
 - Select features of the input
 - Generate neighbors of x by “perturbing” features

Mask/hide words
of a sentence Add noise to
pixels in a region

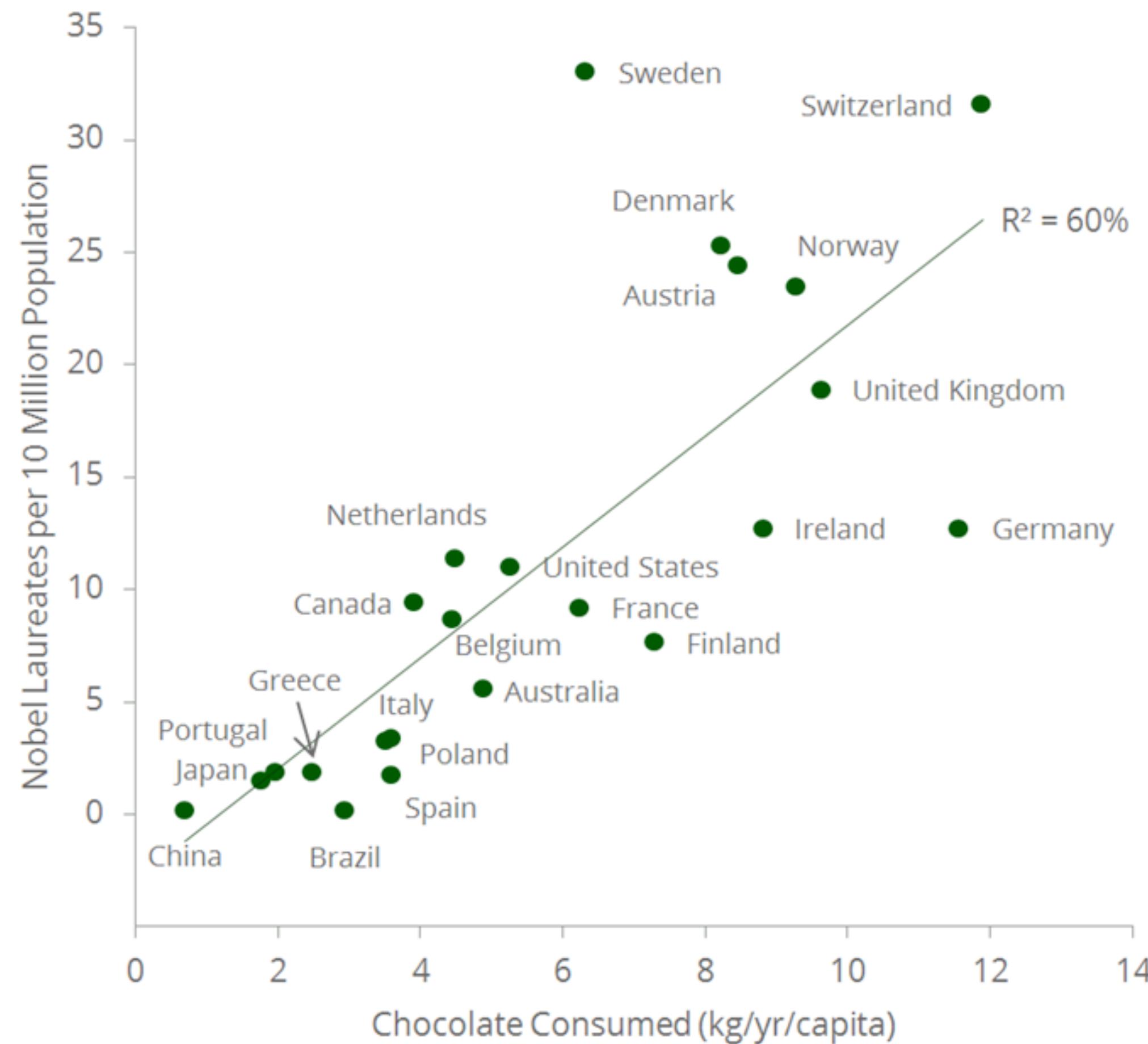
LIME

- For given input x ,
 - Select features of the input
 - Generate neighbors of x by “perturbing” features
 - Train “interpretable model” (e.g., decision tree, linear model) to mimic the output of the original model in the generated neighbors

Correlation

Correlation is not causation

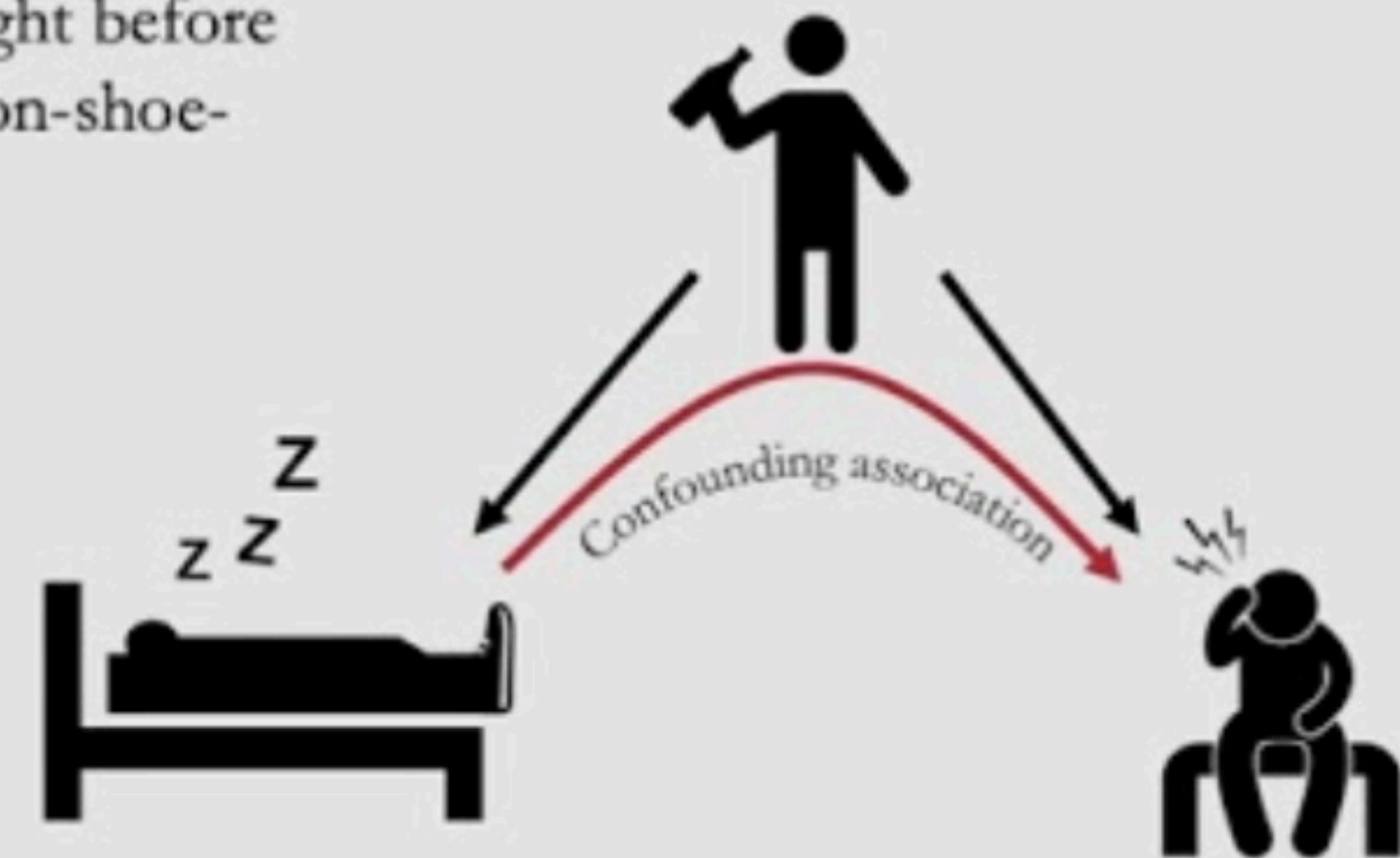
99% of models used (linear or non-linear) make decision based on correlations



Sleeping with shoes on is strongly correlated with waking up with a headache

Common cause: drinking the night before

1. Shoe-sleepers differ from non-shoe-sleepers in a key way
2. Confounding





I USED TO THINK
CORRELATION IMPLIED
CAUSATION.



THEN I TOOK A
STATISTICS CLASS.
NOW I DON'T.

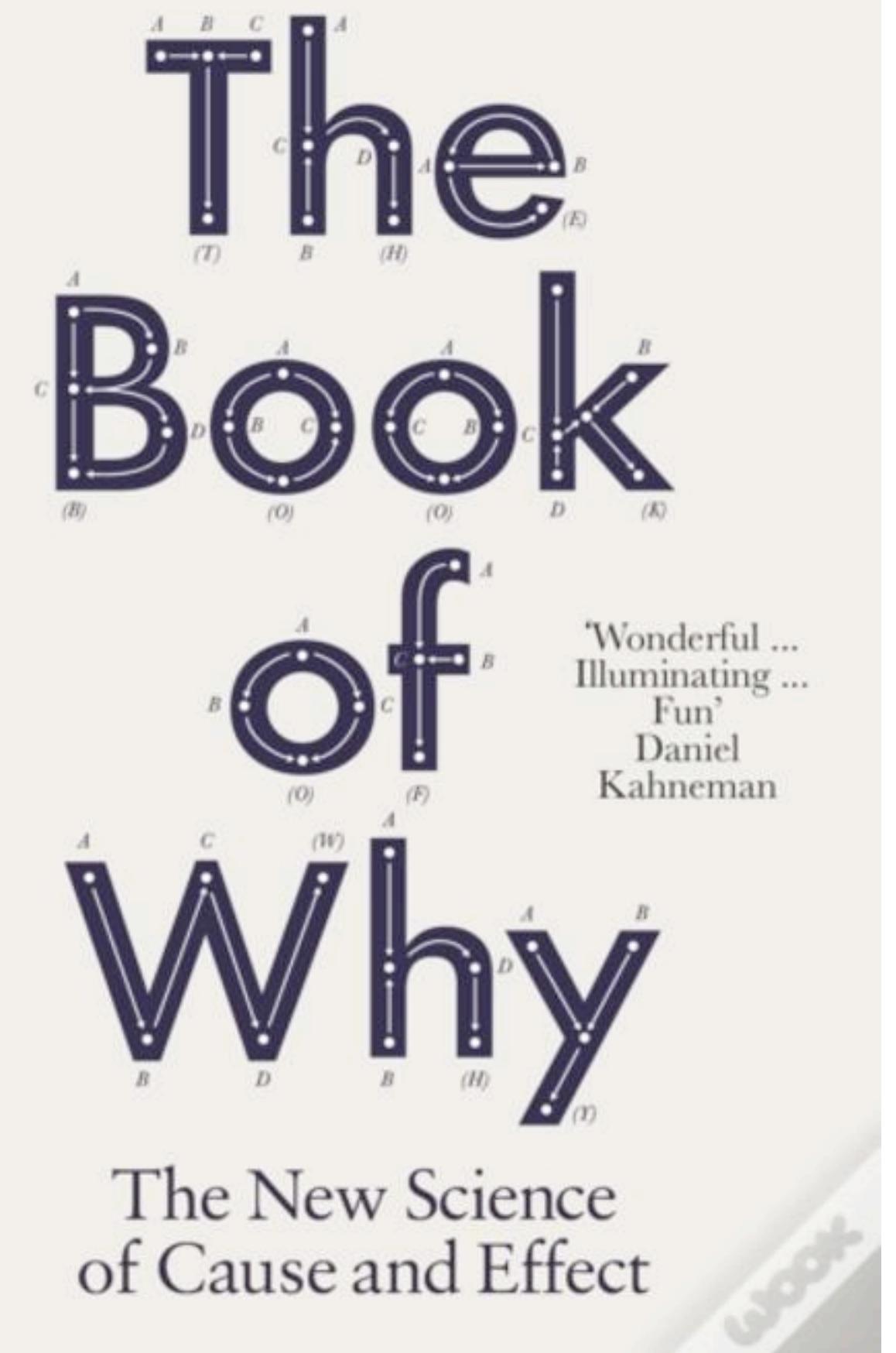


SOUNDS LIKE THE
CLASS HELPED.

| WELL, MAYBE.



Judea Pearl
& Dana Mackenzie

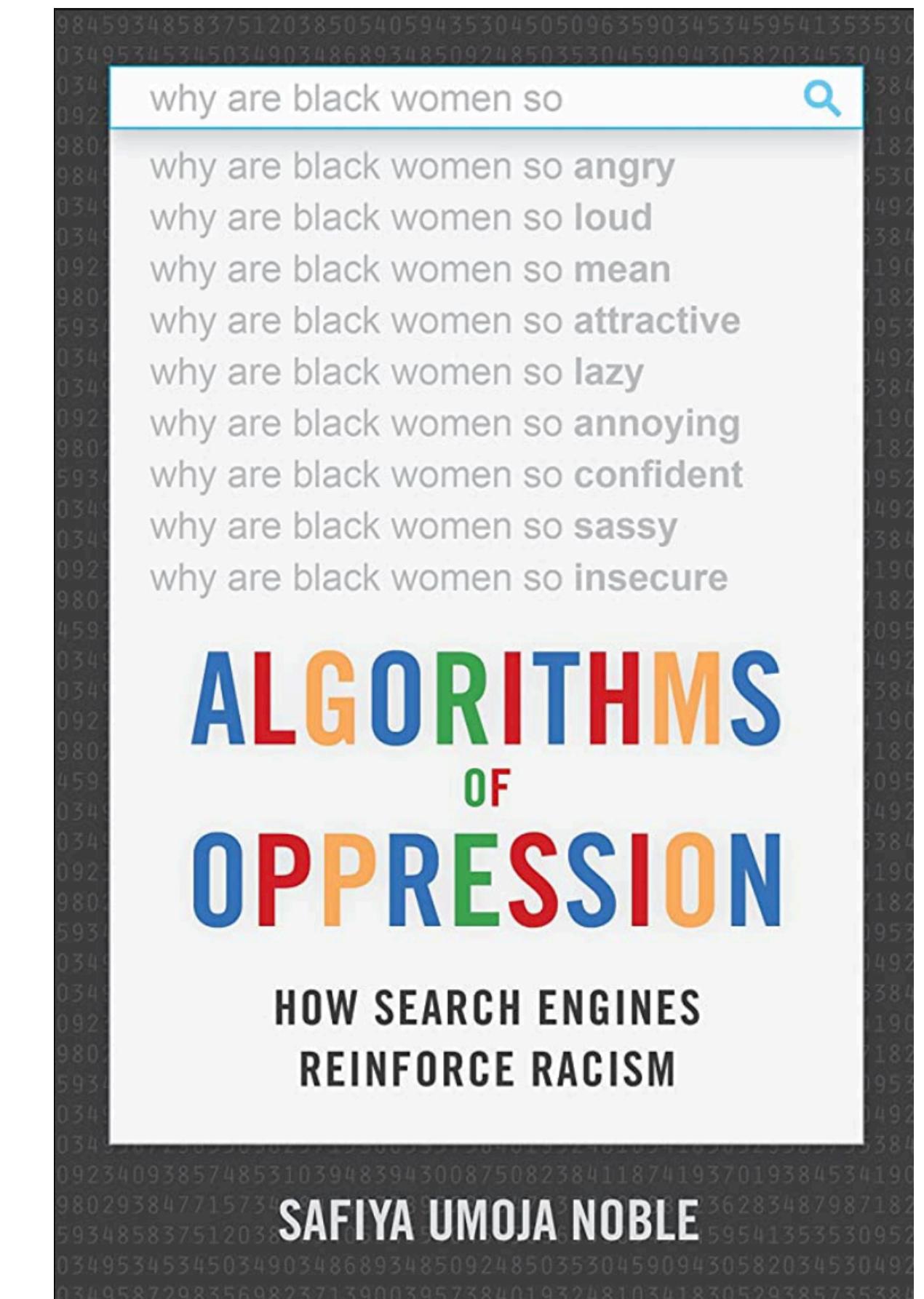
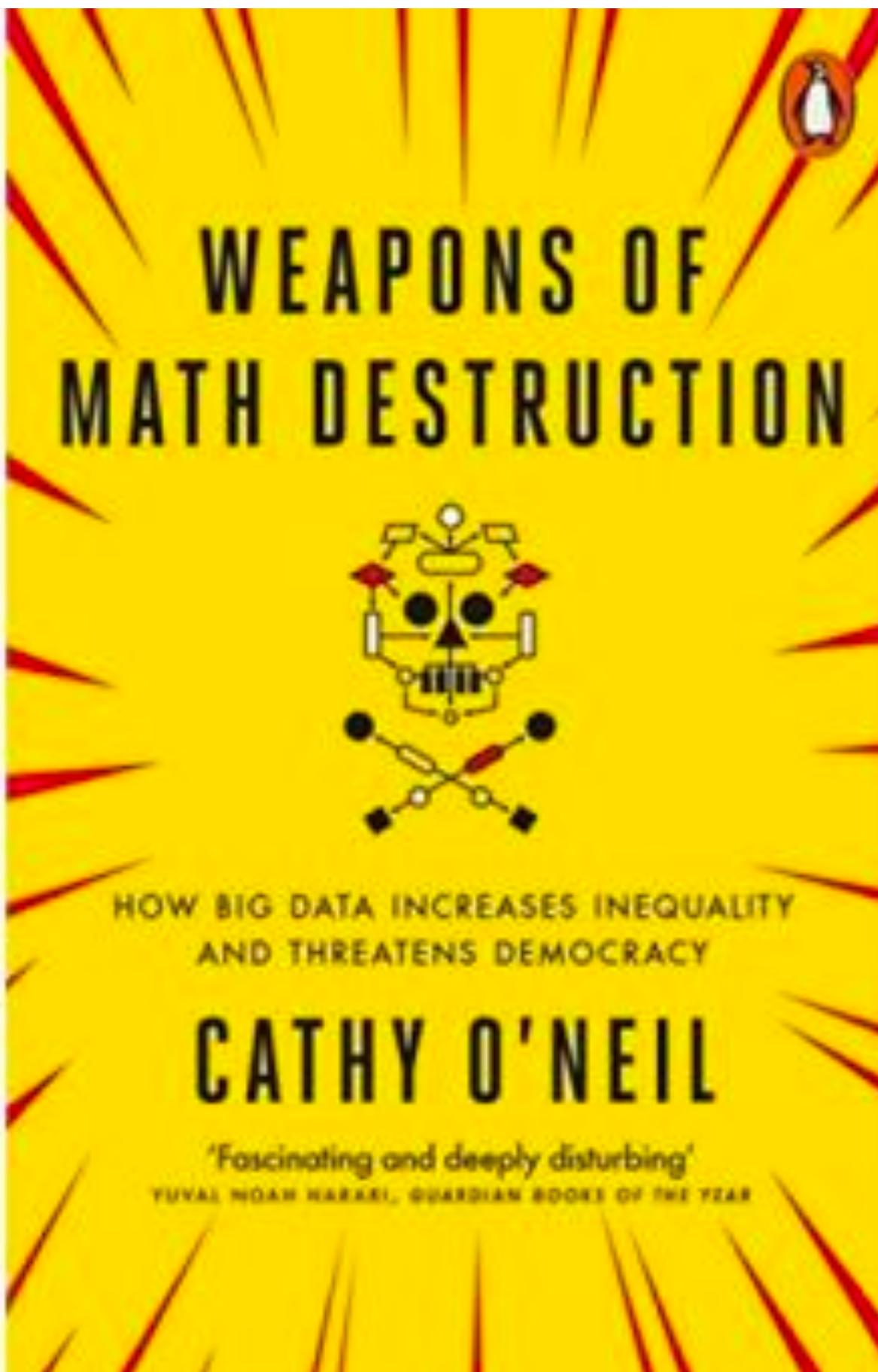


The New Science
of Cause and Effect

Ethics and fairness

Ethics

Beyond performance, trust, and explainability



Ethics

Beyond performance, trust, and explainability

- The extensive use of deep learning models in different aspects of our daily lives raises important ethical concerns
 - Discrimination bias
 - Accessibility
 - Privacy compromise
 - Sustainability

Bias

- What is bias?
- Statistical bias is a systematic tendency which induces differences between results and facts
- Inductive bias (in machine learning) is the set of assumptions that an algorithm uses to make predictions about inputs it has never encountered

Bias

- What is bias?
- **Representation bias** (aka sample bias) is when the training dataset does not include a balanced representation of instances appearing in the test set.
- **Confounding bias** is when there is a distortion of the relation between independent and dependent variables due to a third variable that is independently related to both

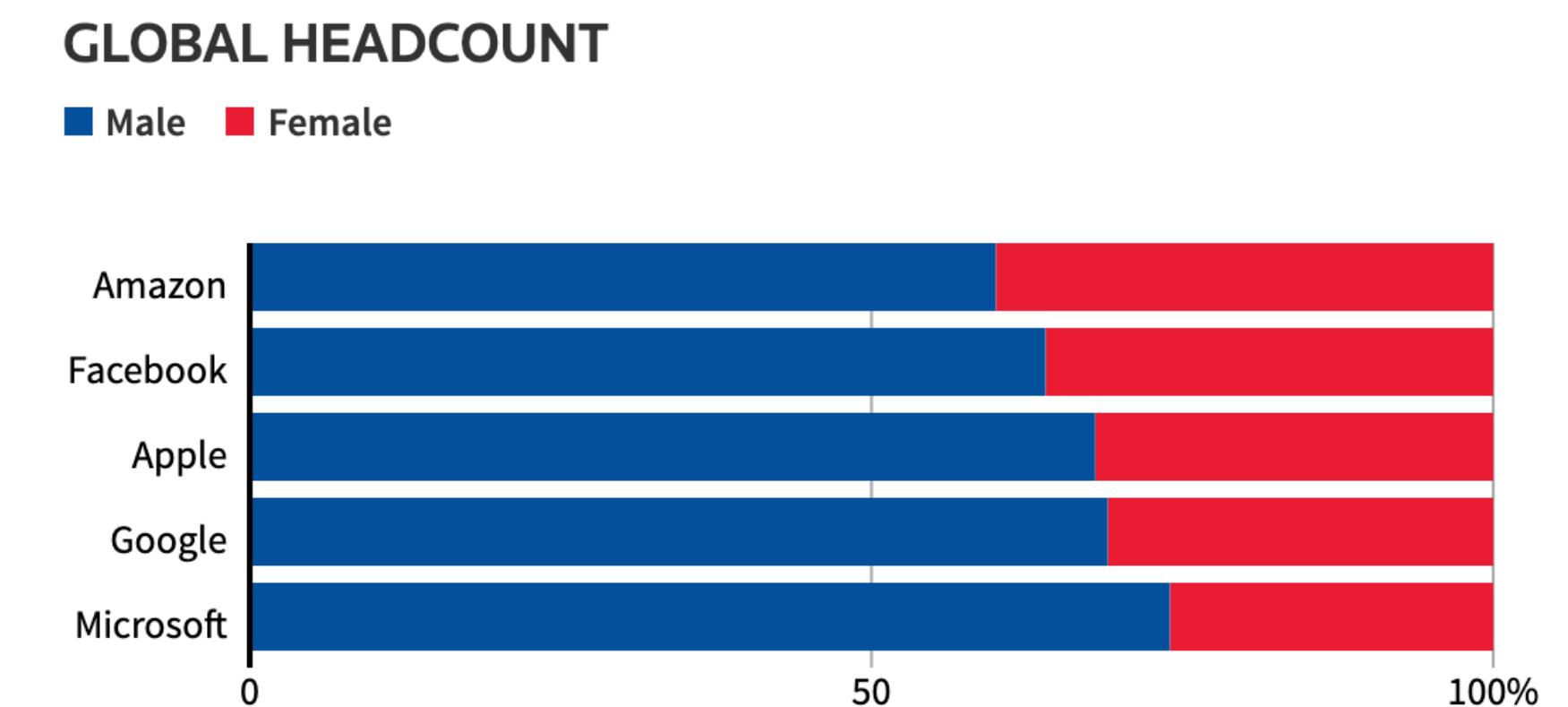
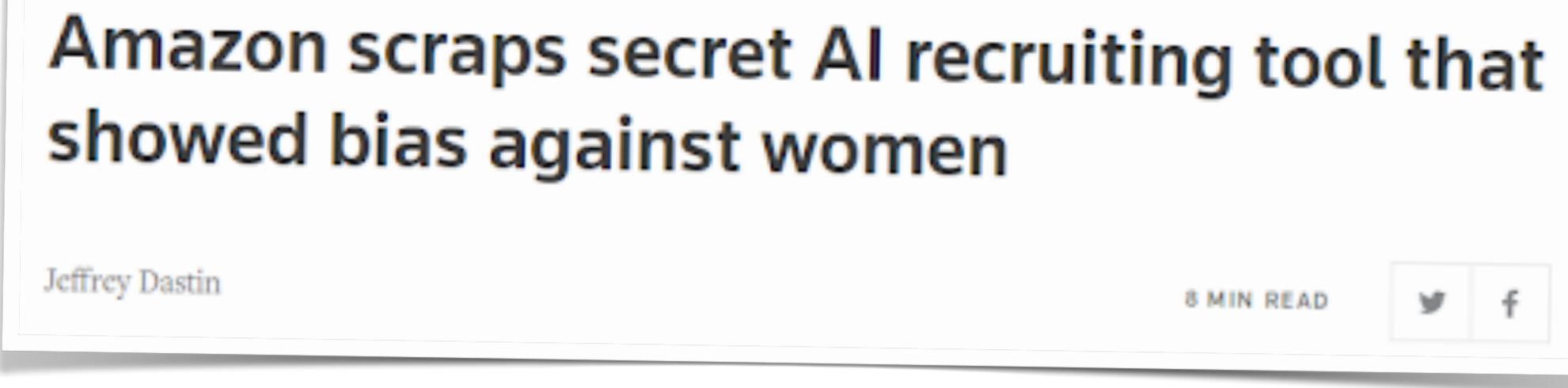
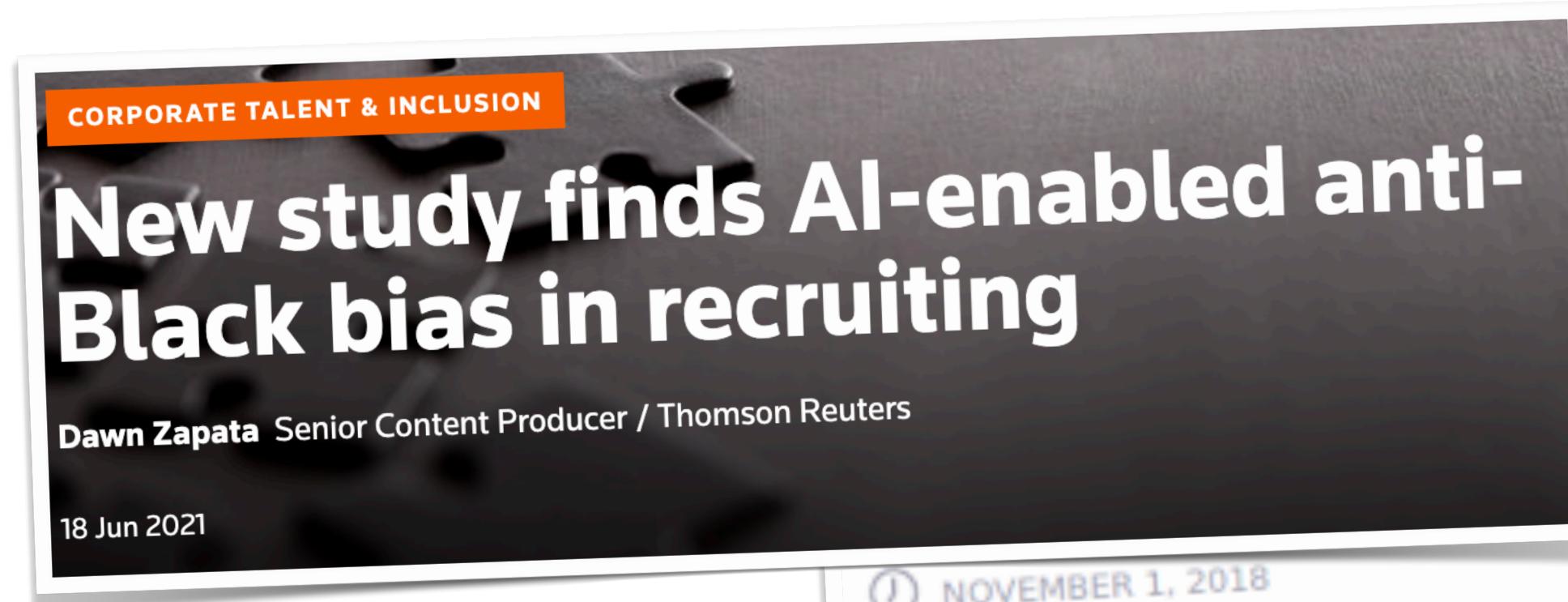
Bias

- What is bias?
- Association bias is when spurious implicit associations occur in the training data but not in the test data
- Observer bias (aka confirmation bias) is when subjective judgements influence the process of data-handling (we “see only what we expect to see”)
- Exclusion bias is when data is removed for being “unimportant”

Bias

- What is bias?
- **Discrimination bias** (aka demographic bias or racial bias) is when the data is skewed to favor a particular demographic group

Discrimination bias



Discrimination bias

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016



Discrimination bias



Snapchat filters don't
work on dark skin



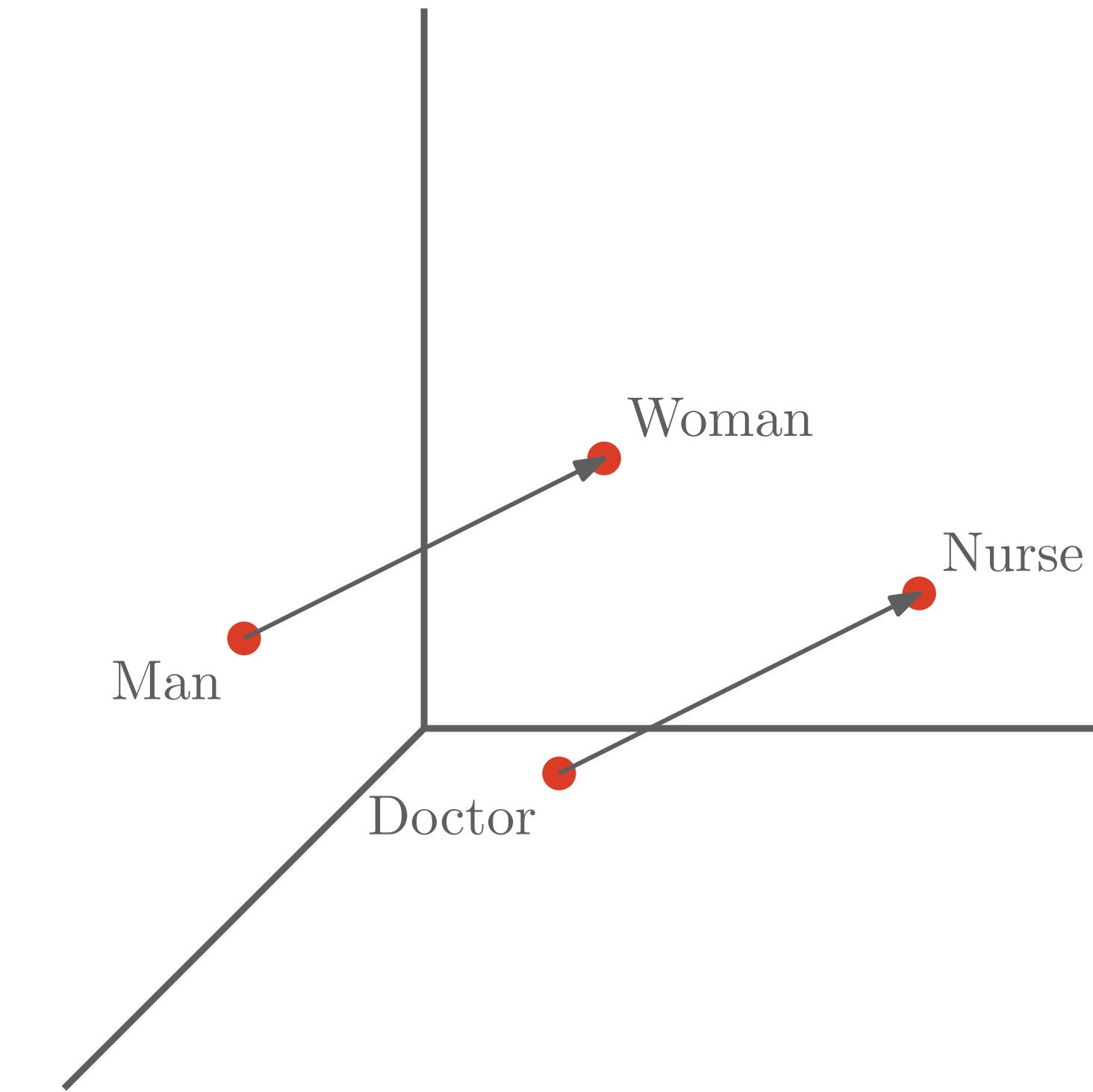
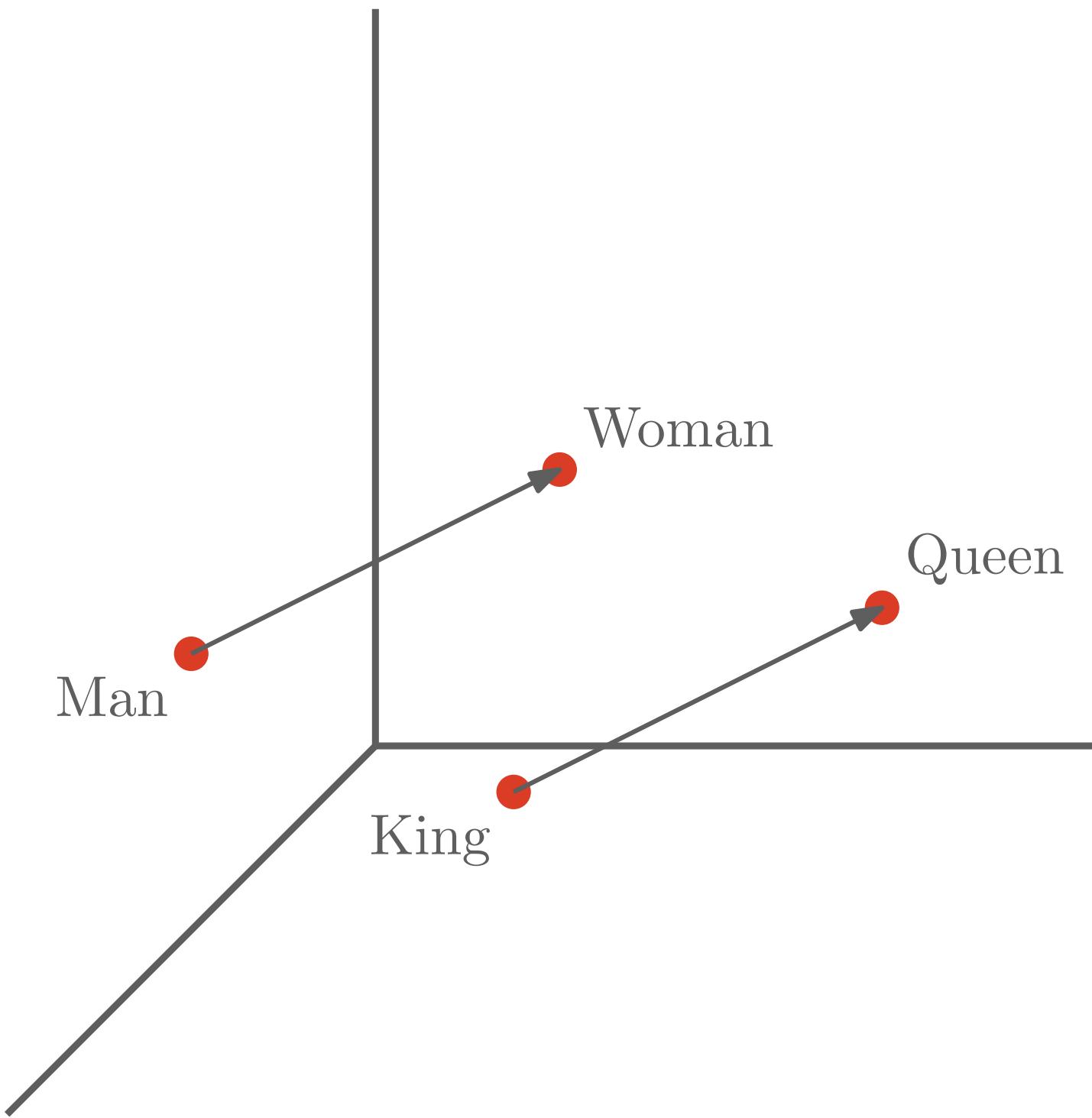
Soap dispensers don't
work on dark skin

Discrimination bias

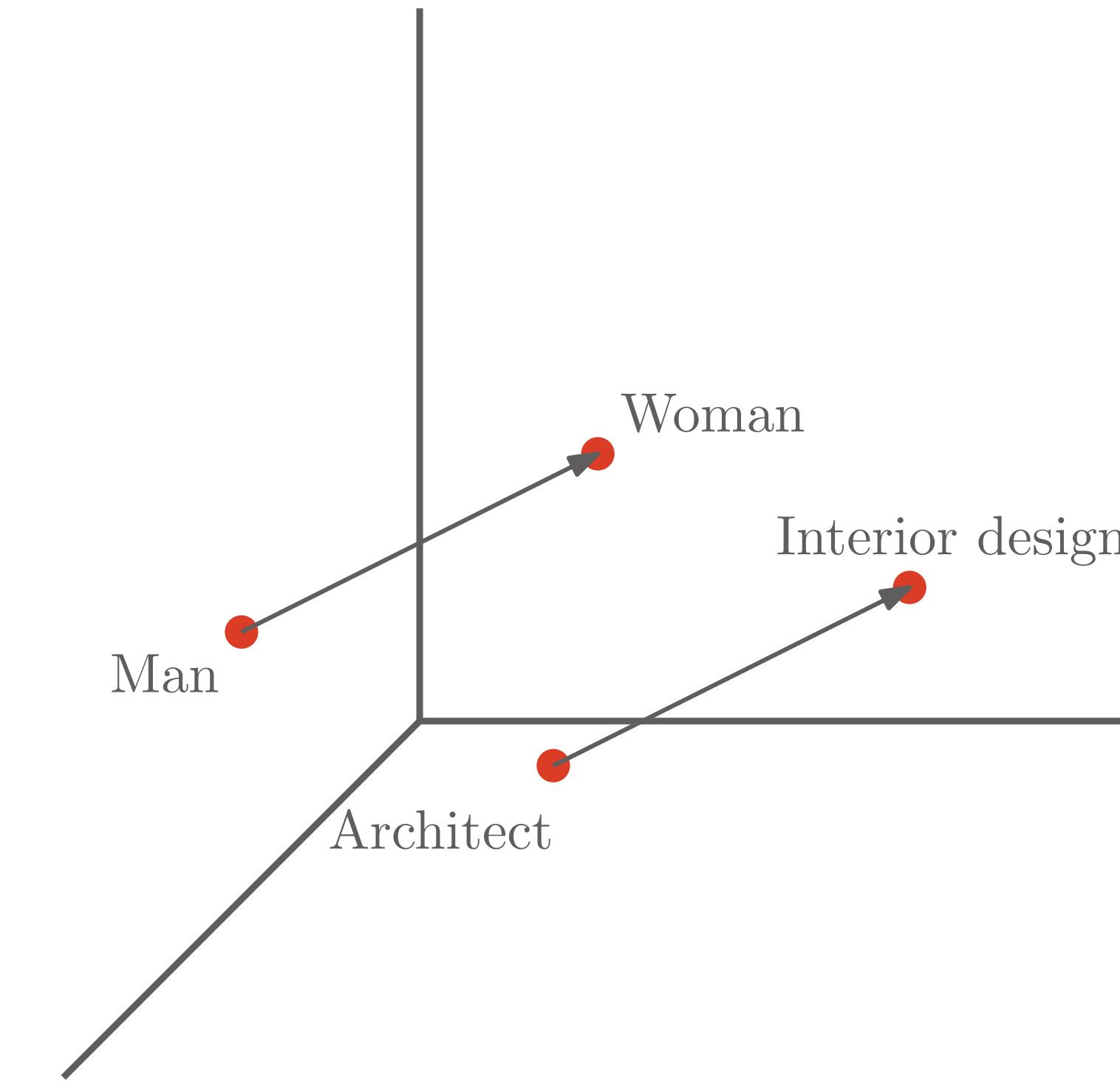
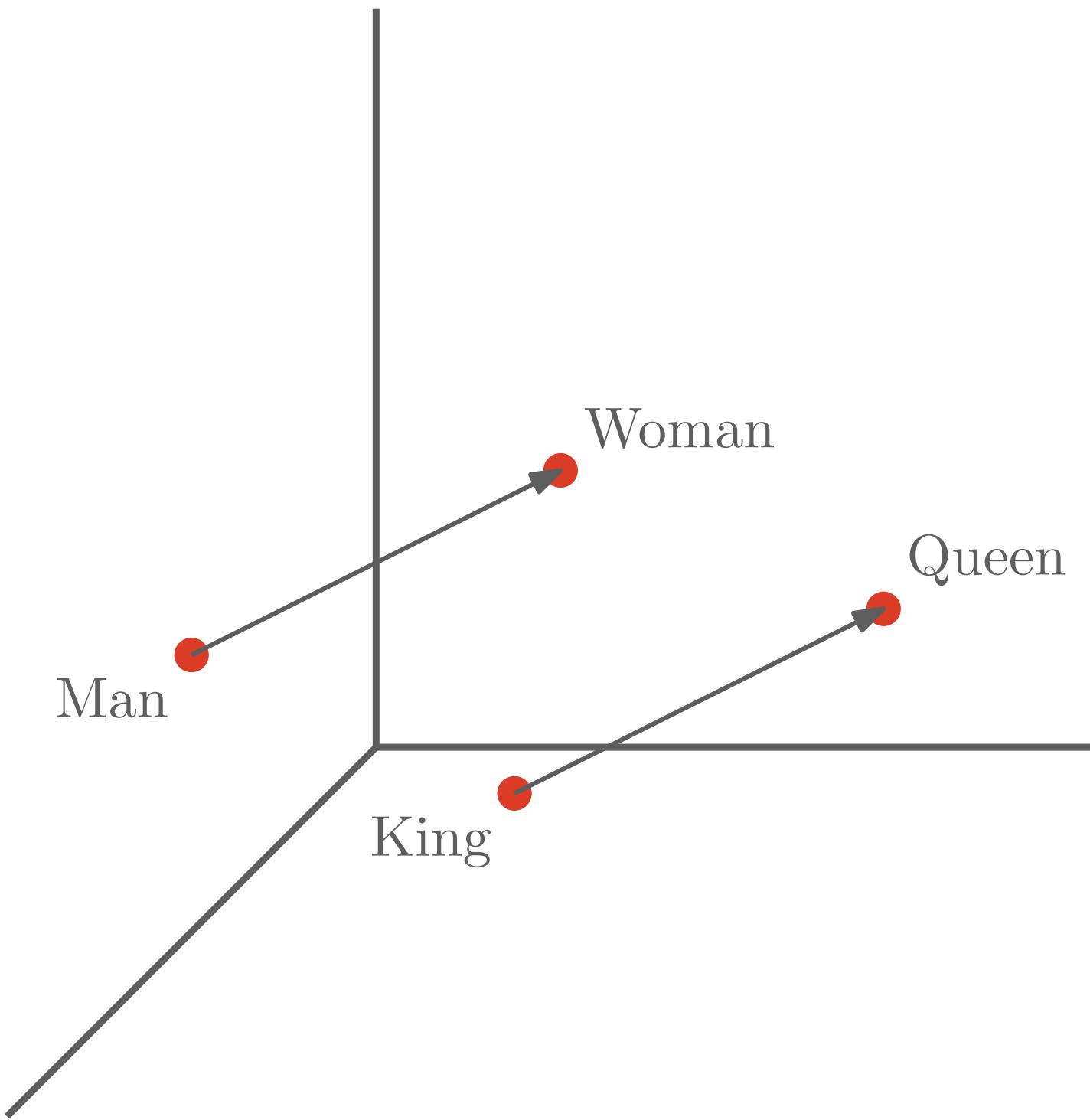
The screenshot shows a translation interface with two panels. The left panel has "DETECT LANGUAGE" and "ENGLISH" selected at the top. The input text is "The nurse had to leave because her shift ended". Below the text are microphone and speaker icons, and a progress bar showing "46 / 5,000". The right panel has "PORTUGUESE" selected at the top. The output text is "A enfermeira teve que sair porque seu turno terminou". Below the text are a speaker icon and sharing icons.

This screenshot is similar to the one above, showing the same translation interface. The input text is "The nurse had to leave because his shift ended". The output text is "A enfermeira teve que sair porque seu turno terminou". A red box highlights the word "enfermeira", and a red arrow points from it to the word "enfermeiro" in the original English input. The UI elements are identical to the first screenshot.

Biased language models



Biased language models



How do we correct biases?

Pre-training:

- Use diverse datasets
- Use social context during annotation
- ...

How do we correct biases?

Training:

Assume feature vectors can be decomposed as

$$\phi(x) = \phi_{\perp}(x) + \phi_g(x)$$

Gender
information



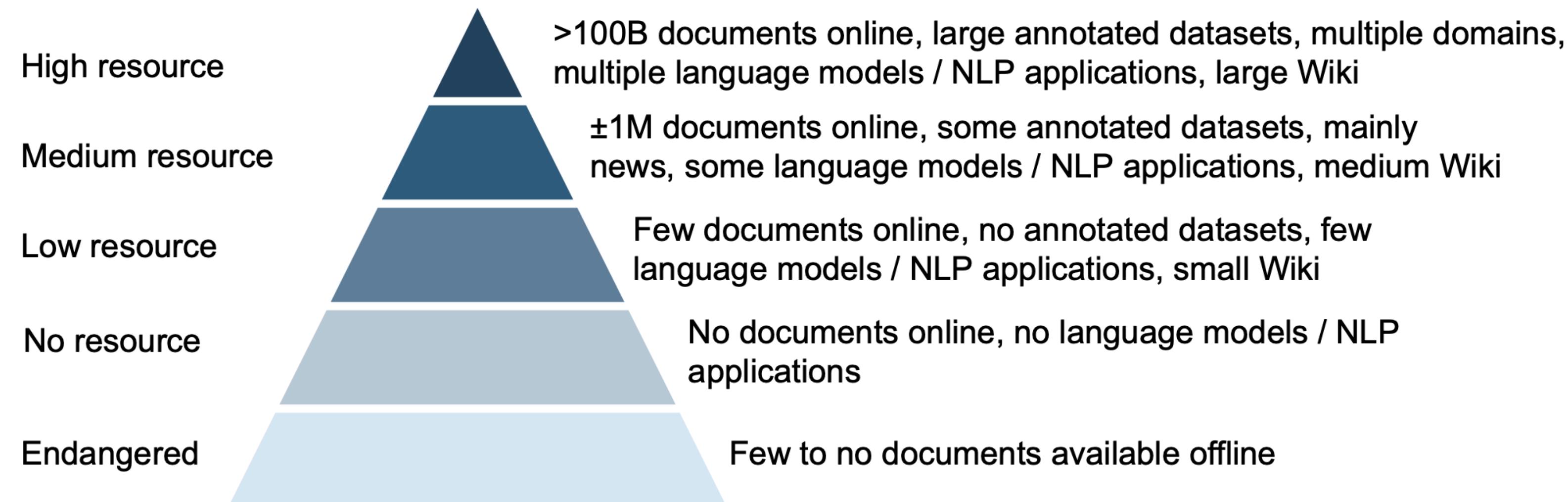
- Loss terms that restrict use of gender information
- Regularization terms that penalize projection in gender dimensions

Accessibility

- Who has access to ML-driven applications?
- Who are ML-driven applications designed for?

An example: NLP

- Are NLP models available for all languages?
- Are multilingual models (e.g., machine translation models) equally efficient/effective for all languages?



Privacy

When is it a concern?

- Should we be tracked?
- Should we be tracked if we consent?
- Even if we consent, how informed is the consent? Do we know how data will be used?
- What if we change our mind?

Privacy

Mumbai-based Agency Leaked Data Of Over 49 Million Instagram Influencers And

The New York Times

Cambridge Analytica and Facebook: The Scandal and the Fallout So Far

Revelations that digital consultants to the Trump campaign misused the data of millions of Facebook users set off a furor on both sides of the Atlantic. This is how The Times covered it.

CPO MAGAZINE

HOME NEWS INSIGHTS RESOURCES

f t in

DATA PRIVACY NEWS · 4 MIN READ

Big Tech Isn't Breaking Any Privacy Rules if There Aren't Rules to Break

DIMITRI SHElest · DECEMBER 27, 2021

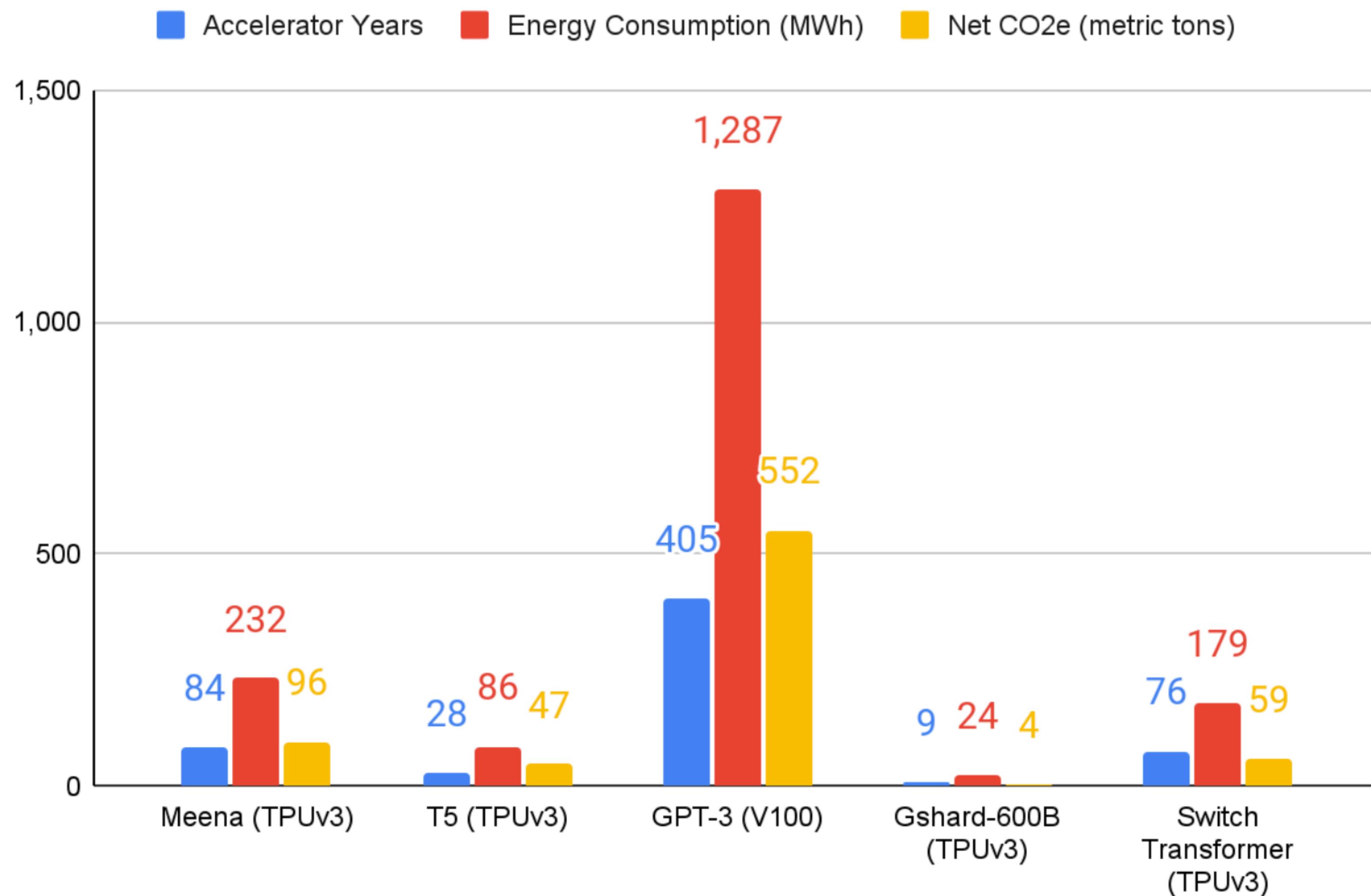
Sustainability

- Big AI and deep learning models consume considerable resources
 - Human resources (engineering, data curation, annotation, etc.)
 - Material resources (hardware)
 - Power consumption
 - CO2 emissions

Training BERT:
Trans-american jet flight



Sustainability



Sustainable AI

- Consider sustainability at every step of the AI/DL life cycle
 - Keep humans in the loop
 - AI can contribute to the development of sustainable solutions towards several of the UN Sustainable Development Goals



Summary

- Generative models are useful to model high-dimensional data
- Both VAEs and GANs fall into the family of differentiable generator networks
- VAEs are trained by maximizing the ELBO
- GANs are trained adversarially

Summary

- The generalization of deep learning technology poses several challenges:
 - How can we develop trustworthy/interpretable models?
 - How can we ensure that deep learning technology development is...
 - ... unbiased?
 - ... accessible?
 - ... sustainable?

To conclude...

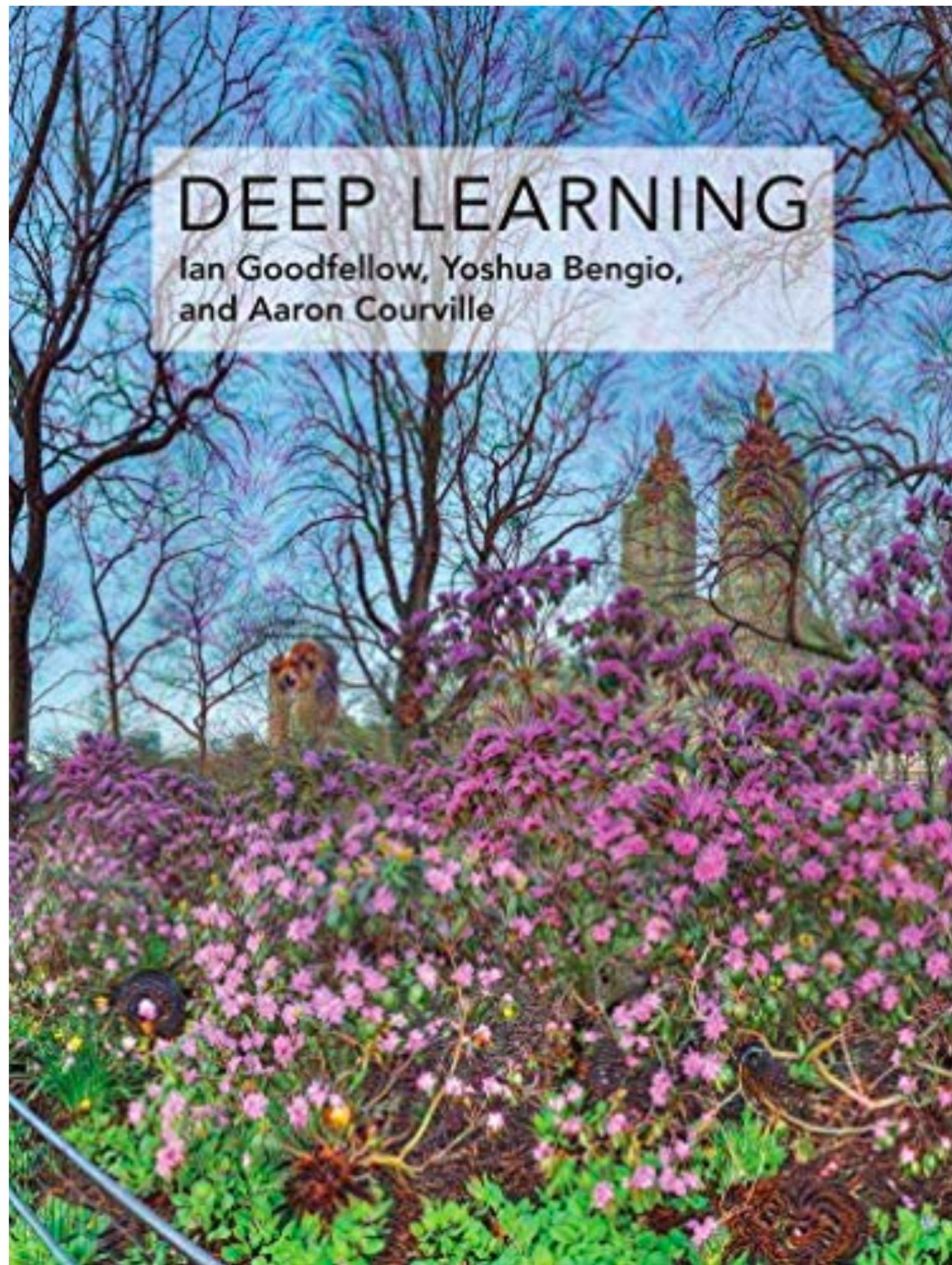
- In the course we covered:
 - Feed-forward neural networks
 - Convolutional neural networks
 - Recurrent neural networks
 - Sequence-to-sequence models and attention mechanisms
 - Generative models

Acknowledgements

A lot of the material used in the course is based on the slides from the MSc course on Deep Learning at Instituto Superior Técnico, authored by André Martins, Mário Figueiredo, and Chrysoula Zerva



References



References

- Bach, S., Binder, A., Montavon, G., Klauschen, F., Muller, K., and Samek, W. “On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation.” *PloS One*, 10(7):e0130140, 2015.
- Bordia, S. and Bowman, S. “Identifying and reducing gender bias in word-level language models.” *CoRR*, abs/1904.03035, 2019.
- Caliskan, A., Bryson, J., and Narayanan, A. “Semantics derived automatically from language corpora contain human-like biases.” *Science*, 356(6334):183-186, 2017.
- Gal, Y., and Z. Ghahramani. “Dropout as a Bayesian approximation: Representing model uncertainty in deep learning.” In *Proc. 33rd Int. Conf. Machine Learning*, pp. 1050-1059, 2016.
- Gonen, H. and Goldberg, Y. “Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them.” *CoRR*, abs/1903.03862, 2019.

References

- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. “Generative adversarial nets.” In *Adv. Neural Information Processing Systems 27*, pp. 2672-2680, 2014.
- Kingma, D., and Welling, M. “Auto-encoding variational Bayes.” *CoRR*, abs/1312.6114, 2013.
- Lakshminarayanan, B., Pritzel, A., and Blundell, C. “Simple and scalable predictive uncertainty estimation using deep ensembles.” In *Adv. Neural Information Processing Systems 30*, pp. 6405-6416, 2017.
- Ribeiro, M., Singh, S., and Guestrin, C. “Why should I trust you? Explaining the predictions of any classifier.” In *Proc. 22nd ACM-SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, pp. 1135-1144, 2016.
- Selvaraju, R., Das, A., Vedantam, R., Cogswell, M., Parikh, D., and Batra, D. “Grad-CAM: Why did you say that?” *CoRR*, abs/1611.07450, 2016.
- Vashisht, S., Upadhyay, S., Tomar, G., and Faruqui, M. “Attention interpretability across NLP tasks.” *CoRR*, abs/1909.11218, 2019.